



**HAL**  
open science

# Generative Supervised Classification Using Dirichlet Process Priors.

Manuel Davy, Jean-Yves Tournet

► **To cite this version:**

Manuel Davy, Jean-Yves Tournet. Generative Supervised Classification Using Dirichlet Process Priors.. IEEE Transactions on Pattern Analysis and Machine Intelligence, 2010, 3 (10), pp.1781-1794. 10.1109/TPAMI.2010.21 . hal-03556758

**HAL Id: hal-03556758**

**<https://hal.science/hal-03556758>**

Submitted on 4 Feb 2022

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



## Open Archive Toulouse Archive Ouverte (OATAO)

OATAO is an open access repository that collects the work of Toulouse researchers and makes it freely available over the web where possible.

This is an author-deposited version published in: <http://oatao.univ-toulouse.fr/>  
Eprints ID: 4150

**To link to this article:** DOI:10.1109/TPAMI.2010.21  
URL: <http://dx.doi.org/10.1109/TPAMI.2010.21>

**To cite this version:** Davy, Manuel and Tourneret, Jean-Yves (2010) *Generative Supervised Classification Using Dirichlet Process Priors*. IEEE Transactions on Pattern Analysis and Machine Intelligence, vol. 32 (n° 10 ). pp. 1781-1794 . ISSN 0162-8828

Any correspondence concerning this service should be sent to the repository administrator: [staff-oatao@inp-toulouse.fr](mailto:staff-oatao@inp-toulouse.fr)

# Generative Supervised Classification Using Dirichlet Process Priors

Manuel Davy, *Member, IEEE*, and Jean-Yves Tournet, *Senior Member, IEEE*

**Abstract**—Choosing the appropriate parameter prior distributions associated to a given Bayesian model is a challenging problem. Conjugate priors can be selected for simplicity motivations. However, conjugate priors can be too restrictive to accurately model the available prior information. This paper studies a new generative supervised classifier which assumes that the parameter prior distributions conditioned on each class are mixtures of Dirichlet processes. The motivations for using mixtures of Dirichlet processes is their known ability to model accurately a large class of probability distributions. A Monte Carlo method allowing one to sample according to the resulting class-conditional posterior distributions is then studied. The parameters appearing in the class-conditional densities can then be estimated using these generated samples (following Bayesian learning). The proposed supervised classifier is applied to the classification of altimetric waveforms backscattered from different surfaces (oceans, ices, forests, and deserts). This classification is a first step before developing tools allowing for the extraction of useful geophysical information from altimetric waveforms backscattered from nonoceanic surfaces.

**Index Terms**—Supervised classification, Bayesian inference, Gibbs sampler, Dirichlet processes, altimetric signals.

## 1 INTRODUCTION

STATISTICAL classification has received considerable attention in the literature (see, for instance, [1], [2], [3], [4], [5], and references therein). Among all possible classification strategies, the Bayes decision rule is known to provide a minimal probability of error (or, more generally, a minimal expected value of an appropriate loss function). However, the Bayes classifier assumes knowledge of the class-conditional densities and the prior class probabilities. When the class-conditional densities are unknown, a classifier is generally built by using training patterns from the different classes following supervised learning. This paper focuses on supervised parametric classification. Each feature vector belonging to one of  $J$  identified classes is assumed to be the random outcome of some stochastic process tuned by an unknown parameter vector. In this setting, training the classifier consists of learning the parameter class-conditional posterior distribution (see Section 2) so as to compute Bayes factors [6]. We focus in this paper on generative classifiers that provide better performance than discriminative classifiers for small training set sizes [7].

Learning a parameter posterior distribution is generally a difficult problem. The parameter prior distributions have to be sufficiently general to model accurately the parameter class-conditional distributions without yielding untractable computational complexity. As a consequence, advanced learning strategies have to be developed to handle the

possible multiple maxima of the objective function. These difficulties explain why fully Bayesian generative supervised classifiers have received less attention than other classifiers. A simple generative training strategy known as naive Bayes classification assumes that all features are independent inside a given class. This strategy has been used successfully in many practical applications, including text classification [8] and object detection [9]. Semi-naive Bayes classifiers decompose the input features into statistically independent subsets with statistical dependency within each subset [10]. This paper studies a new generative classifier based on Dirichlet process priors. These priors are known to enable the modeling of a large class of probability density functions (pdfs) (e.g., see [11], [12], [13], [14], [15], and references therein). The naive Bayes and semi-naive Bayes classifiers mentioned above are special cases of the proposed generative classifier based on Dirichlet process priors.

### 1.1 Contributions

A generative classifier based on hierarchical Bayesian models and Markov chain Monte Carlo (MCMC) methods has been recently studied in [16]. The Bayesian approach of [16] is adapted to cases where a precise generative data model can be defined. However, the classifier may be difficult to generalize to any classification problem and it can lead to high-computational complexity when conjugate priors cannot be used for the unknown parameters. The present paper studies an alternative supervised generative classifier involving Dirichlet process priors. Our approach assumes a nonparametric flexible parameter prior distribution with controlled complexity, known as the Dirichlet Process Mixture (DPM) [17]. Our main contribution is to derive DPM-based models for generative supervised classification as well as efficient MCMC algorithms, allowing one to estimate the unknown parameters of these models. It is interesting to mention here the work by Shahbaba and Neal [18], which is also related to supervised classification using DPMs. However, there are major

- M. Davy is with VEKIA, 165 Avenue de Bretagne, 59000 Lille, France. E-mail: mdavy@vekia.fr.
- J.-Y. Tournet is with the University of Toulouse, ENSEEIHT-IRIT-TéSA, 2 Rue Camichel, BP 7122, 31071 Toulouse Cedex 7, France. E-mail: jean-yves.tournet@enseeiht.fr.

differences between [18] and the proposed approach, as summarized below:

- Shahbaba and Neal have addressed classification in a discriminative context without using any generative model. Conversely, the proposed classification strategy studied here is generative allowing us to efficiently use the data and a parametric model associated to these data.
- Instead of modeling the parameter prior by a DPM, Shahbaba and Neal propose modeling the joint distribution of the training data and the class labels by a DPM.
- Shababa and Neal assume that the training data have a Gaussian distribution. This assumption is not necessary for the proposed algorithm. For instance, the radar data studied in the simulation section are distributed according to gamma distributions.

## 1.2 Paper Organization

A comprehensive glossary of the main notations is provided in Appendix A. Section 1 introduces the problem addressed in this paper. Section 2 presents the proposed parametric hierarchical Bayesian model used for supervised classification. Section 3 recalls some useful properties of Dirichlet processes. Section 4 studies original prior models for supervised classification. These models assume that the parameter priors for each class are mixtures of Dirichlet processes. Consequently, these models do not assume that the conditional class densities have a known parametric shape as is often the case in classification problems (see [2, p. 84] for details). Section 5 applies the proposed Bayesian methodology to the classification of synthetic and real data. Conclusions are reported in Section 6.

## 2 BAYESIAN SUPERVISED CLASSIFICATION

This section recalls the principle of Bayesian supervised classification, using a parametric hierarchical model.

### 2.1 Problem Settings

Assume we want to classify (vectorial) data into  $J$  classes denoted as  $C_1, \dots, C_J$ . We assume further that a training set is available for each class  $C_j$ ,  $j = 1, \dots, J$ , denoted as  $\mathbf{X}_j = \{\mathbf{x}_{1,j}, \dots, \mathbf{x}_{N_j,j}\}$ . Inside each class, we assume the following generative model:

$$\mathbf{x}_{i,j} = M_j(\theta_{i,j}, \epsilon_{i,j}), \quad (1)$$

where  $M_j(\cdot)$  is a known function, the parameters  $\theta_{i,j}$  live in a space  $\Theta_j \subset \mathbb{R}^d$  ( $d$  is the dimension of the parameter space),  $i = 1, \dots, N_j$ , and  $\epsilon_{i,j}$  is some stochastic noise. Note that the noise needs not be additive, e.g., the case of a multiplicative speckle noise will be considered in Section 5. In probabilistic terms,

$$\mathbf{x}_{i,j} \sim p_j(\mathbf{x}_{i,j}|\theta_{i,j}), \quad (2)$$

where  $p_j(\mathbf{x}_{i,j}|\theta_{i,j})$  is the pdf of  $\mathbf{x}_{i,j}$  (belonging to the class  $C_j$ ) conditioned upon  $\theta_{i,j}$ .

The Bayesian algorithms, derived in this paper, estimate the posterior distribution of the parameter vector  $\theta_j$  from the training set  $\mathbf{X}_j$  for each class  $C_j$ ,  $j = 1, \dots, J$ . The parameter

posterior distribution for class  $C_j$  will be denoted as  $p(\theta_j|\mathbf{X}_j)$  and its estimation, detailed in the next section, will be referred to as Bayesian learning.

### 2.2 Parametric Hierarchical Approach

Within the parametric hierarchical Bayesian settings, learning  $p(\theta_j|\mathbf{X}_j)$  requires the definition of a prior distribution over the parameter space  $\Theta_j$ , denoted as  $p_j(\theta_j|\phi_j)$ , with hyperparameter  $\phi_j$ . The posterior distribution of  $\theta_j$  can then be classically expressed as:<sup>1</sup>

$$p(\theta_j|\mathbf{X}_j) = \int p_j(\theta_j|\phi_j)p(\phi_j|\mathbf{X}_j)d\phi_j, \quad (3)$$

where

$$\begin{aligned} p(\phi_j|\mathbf{X}_j) &= \int \dots \int p(\theta_{1,j}, \dots, \theta_{N_j,j}, \phi_j|\mathbf{X}_j)d\theta_{1,j} \dots d\theta_{N_j,j} \\ &\propto \int \dots \int p_j(\phi_j) \prod_{i=1}^{N_j} [p_j(\mathbf{x}_{i,j}|\theta_{i,j})p_j(\theta_{i,j}|\phi_j)] \\ &\quad d\theta_{1,j} \dots d\theta_{N_j,j} \end{aligned} \quad (4)$$

by Bayes rule ( $\propto$  means ‘‘proportional to’’). As can be seen, a prior distribution over the hyperparameter, denoted as  $p_j(\phi_j)$ , has to be defined to determine the posterior parameter distribution  $p(\theta_j|\mathbf{X}_j)$  in (3). Once  $p(\theta_j|\mathbf{X}_j)$  is known, a new observation vector  $\mathbf{x}$  can be classified in view of the predictive densities

$$p(\mathbf{x}|\mathbf{X}_j) = \int p_j(\mathbf{x}|\theta_j)p(\theta_j|\mathbf{X}_j)d\theta_j, \quad (5)$$

for  $j = 1, \dots, J$ . For instance, assuming equal class prior probabilities, the Bayesian maximum a posteriori (MAP) classifier assigns  $\mathbf{x}$  to the class maximizing  $p(\mathbf{x}|\mathbf{X}_j)$ . The Bayesian algorithms derived in this paper estimate the posterior distributions  $p(\theta_j|\mathbf{X}_j)$  from the training set  $\mathbf{X}_j$ .

### 2.3 Bayesian Computations

In this section, we drop the class subscript as the derivations are performed in each class in turn, independently of the other classes. Equivalently, we assume that the training samples in a class  $C_j$  do not provide any information regarding parameter  $\theta_i$  for  $i \neq j$  (see [2, p. 85] for a similar assumption). The classification of a given observation vector requires the evaluation of multidimensional integrals in (3) and (4). These integrals cannot be computed in closed-form in the general case. However, they can be evaluated by using Monte Carlo methods as proposed in [16]. The resulting algorithm is briefly recalled below (the reader is invited to consult [16] for more details).

Assume that a set of samples  $\theta^{(l)}$ ,  $l = 1, \dots, L$ , distributed according to  $p(\theta|\mathbf{X})$  is available. Note that  $p(\theta|\mathbf{X})$  is defined as in (3), where the training set  $\mathbf{X}_j$  has been replaced by  $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_N)$  for notation simplicity (the class subscript has been omitted). Then, the integral in (5) can be estimated as follows:

1. For the sake of notational clarity, we adopt a slight abuse of notation by denoting Lebesgue integrals over distributions as standard integrals. These integrals should be understood as sums whenever the parameter to be integrated is discrete.

$$p(\mathbf{x}|\mathbf{X}) = \int p(\mathbf{x}|\theta)p(\theta|\mathbf{X})d\theta \approx \frac{1}{L} \sum_{l=1}^L p(\mathbf{x}|\tilde{\theta}^{(l)}). \quad (6)$$

A usual way of generating samples  $\tilde{\theta}^{(l)}$  distributed according to  $p(\theta|\mathbf{X})$  is summarized in Algorithm 1.

**Algorithm 1.** Sampling from the class parameter posterior distribution

For  $l = 1, \dots, L$ , do

1) sample  $\tilde{\phi}^{(l)} \sim p(\phi|\mathbf{X})$  given in (4)

2) sample  $\tilde{\theta}^{(l)} \sim g(\theta|\tilde{\phi}^{(l)})$

Algorithm 1 requires being able to sample from  $p(\phi|\mathbf{X})$ . In order to do so, one generates samples distributed according to the joint posterior distribution  $p(\theta, \phi|\mathbf{X})$ , where  $\theta = \{\theta_1, \dots, \theta_N\}$ . The parameters  $\theta_i$ s are then marginalized out by only considering the component  $\tilde{\phi}^{(l)}$  of the full set of sampled parameters and hyperparameters. In order to sample from the full posterior  $p(\theta, \phi|\mathbf{X})$ , one implements the following Gibbs sampler (where MH stands for ‘‘Metropolis-Hastings’’):

**Algorithm 2.** Gibbs sampler for Bayesian Supervised classification

% Step 0: Initialization

• Sample  $\tilde{\phi}^{(0)} \sim h(\phi)$

% step 1: Iterations. For  $l = 1, \dots, L$ , do

- For  $i = 1, \dots, N$ , sample  $\tilde{\theta}_i^{(l)} \sim p(\theta_i|\mathbf{x}_i, \tilde{\phi}^{(l-1)}) \propto p(\mathbf{x}_i|\theta_i)g(\theta_i|\tilde{\phi}^{(l-1)})$  using, e.g., an MH step

- Sample  $\tilde{\phi}^{(l)} \sim p(\phi|\tilde{\theta}^{(l)}) \propto p(\phi) \prod_{i=1}^N g(\tilde{\theta}_i^{(l)}|\phi)$  using, e.g., an MH step

## 2.4 Comments

Though straightforward to implement, Bayesian supervised classification requires the definition of several distributions. In practice, the distribution  $p(\mathbf{x}|\theta)$  (referred to as likelihood) is imposed by empirical considerations. The prior  $p(\theta|\phi)$  is chosen for its conjugacy as well as for physical considerations. However, its choice remains arbitrary in large parts. In the following, we propose a nonparametric parameter prior constructed from DPMs.

## 3 DIRICHLET PROCESSES AND DIRICHLET PROCESS MIXTURES

### 3.1 Dirichlet Processes

Dirichlet processes are characterized by a base distribution  $F_0$  defined on an appropriate measurable space  $(\Theta, \mathcal{B})$  and a positive real number  $\alpha$ . More precisely, a Dirichlet process with base distribution  $F_0$  and concentration parameter  $\alpha$ , denoted as  $F \sim \mathcal{DP}(\alpha, F_0)$ , is the distribution of a random probability measure  $F$  over  $(\Theta, \mathcal{B})$  such that, for Borel sets  $A_1, \dots, A_m$  forming a partition<sup>2</sup> of  $\Theta$ ,

$$[F(A_1), \dots, F(A_m)] \sim \mathcal{D}(\alpha F_0(A_1), \dots, \alpha F_0(A_m)), \quad (7)$$

where  $\mathcal{D}$  is the standard Dirichlet distribution [17], [19]. Equivalently, a Dirichlet process  $F$  can be defined as a discrete distribution

2. A partition of a space  $\Phi$  is a set of subsets  $A_1, \dots, A_m$  with arbitrary  $m$  such that  $A_i \cap A_j = \emptyset$  for any  $i \neq j$  and  $A_1 \cup \dots \cup A_m = \Phi$ .

$$F = \sum_{k=1}^{\infty} \omega_k \delta_{U_k}, \quad (8)$$

where the probabilities  $\omega_k$  are scalar nonnegative values that sum to one and the  $U_k$  are random variables living in the space of the base distribution. The random variables  $(U_1, U_2, \dots)$  are independent of the probabilities  $(\omega_1, \omega_2, \dots)$  and i.i.d. among themselves. They can be generated recursively as follows:  $U_1 \sim F_0$ ,  $\omega_1 \sim \mathcal{B}(\omega_1; 1, \alpha)$ , and, for  $k = 2, 3, \dots$ :

- $U_k \sim F_0(\cdot)$ ,
- $\beta_k \sim \mathcal{B}(\beta_k; 1, \alpha)$  and  $\omega_k = \beta_k \prod_{k'=1}^{k-1} (1 - \beta_{k'})$ ,

where  $\mathcal{B}(\cdot; a, b)$  denotes the beta distribution with parameters  $a$  and  $b$  (defined in [2, p. 109]). This procedure is usually referred to as *stick breaking representation* [11], [12], [20] and explains the role of  $F_0$  and  $\alpha$ :  $F_0$  determines the locations of the discrete components of  $F$ , while  $\alpha$  determines the variance of the  $\omega_k$ s. The reader is invited to consult [21] to understand the role of parameter  $\alpha$ . A first analysis is conducted in [21] with a fixed value of  $\alpha$ , i.e.,  $\alpha = 1$ . The parameter  $\alpha$  is then varied to understand how sensitive the results are to this parameter.

### 3.2 Polya Urn Representation

Consider a set of independent identically distributed (i.i.d) variables  $\{\phi_1, \dots, \phi_N\}$  distributed according to  $F(\cdot)$ , where  $F \sim \mathcal{DP}(\alpha, F_0)$ . Then, for any  $i = 1, \dots, N$ ,

$$p(\phi_i|\phi_{-i}, F_0, \alpha) = \frac{\alpha}{\alpha + N - 1} F_0(\phi_i) + \frac{1}{\alpha + N - 1} \sum_{i'=1, i' \neq i}^N \delta_{\phi_{i'}}(\phi_i), \quad (9)$$

where  $\phi_{-i} = \{\phi_{i'}\}_{i'=1, \dots, N, i' \neq i}$ . This formula is known as the *Polya urn representation*, which also makes clear the role of  $F_0$  and  $\alpha$  [22]. The base distribution  $F_0$  is used to sample the location of clusters. Large values of  $\alpha$  cause many clusters to be created (the probability to sample from  $F_0$  in (9) increases with  $\alpha$ ). Conversely, small values of  $\alpha$  favor sampling from existing clusters and, overall, few clusters are generated. In theory, it is possible to have an infinite number of clusters. However, for a small value of  $\alpha$ , only a few of them have a significant weight  $\omega_k$ , and the probability that at least one  $\phi_i$  belongs to such a cluster becomes extremely small.

The Polya urn representation can also be written for the posterior distribution of  $F$  given  $\{\phi_1, \dots, \phi_N\}$ , as follows:

$$p(F|\phi_1, \dots, \phi_N, F_0, \alpha) = \mathcal{DP}(F_1, \alpha + N), \quad (10)$$

where

$$F_1(\phi) = \frac{\alpha}{\alpha + N} F_0(\phi) + \frac{1}{\alpha + N} \sum_{i=1}^N \delta_{\phi_i}(\phi). \quad (11)$$

### 3.3 Dirichlet Process Mixtures

A DPM is characterized by a DP  $F$  (with base distribution  $F_0$  and concentration parameter  $\alpha$ ) and a mixture distribution  $f$ . Precisely, a random variable  $\theta$  is distributed according to a DPM whenever it is generated as follows:

1. sample  $F \sim \mathcal{DP}(\alpha, F_0)$ ,
2. sample  $\phi \sim F(\cdot)$ , and
3. sample  $\theta \sim f(\theta|\phi)$ .

By iterating Steps 2 and 3, one builds a family of random variables distributed according to a DPM with mixture distribution  $f$ , base distribution  $F_0$ , and parameter  $\alpha$ . The role of  $f$  is to smooth the DP and to turn it into a continuous support distribution (whereas the DP  $F$  has discrete support). Each iteration  $i$  ( $i = 1, \dots, N$ ) of Step 2 yields a parameter  $\phi_i$  that determines the clusters  $\theta_i$ . The Polya urn representation (9) shows that several  $\phi_i$ s can share the same location. This can also be seen in the stick-breaking representation (8) emphasizing the locations  $U_{ks}$  and their probabilities  $\omega_k$ .

Similarly to the finite mixture case studied in [23], it is possible to introduce a latent variable  $z_i$  to each  $\theta_i$  ( $i = 1, \dots, N$ ). This consists of replacing Step 2 above by:

1. Sample  $z \sim p(z|F)$ , where  $p(z = k|F) = \omega_k$  and where  $\omega_k$  comes from the stick-breaking representation of  $F$ .
2. Set  $\phi = U_z$ , where  $U_z$  comes from the stick-breaking representation of  $F$ .

Finally, the following alternative representation of a DPM measure  $G(\theta)$  can be useful:

$$G(\theta) = \int_{\Theta} F(\theta|\phi) dF(\phi) \text{ with } F \sim \mathcal{DP}(F_0, \alpha), \quad (12)$$

where  $F(\theta|\phi)$  is the cumulative distribution function associated to the density  $f(\theta|\phi)$ .

#### 4 USING DPMS TO MODEL THE PARAMETER PRIOR

This section introduces an original Bayesian supervised classification classifier based on DPM priors, referred to as DPM classifier. DPM classification consists of replacing the parameter prior distribution  $\int_{\Phi} g(\theta|\phi)h(\phi)d\phi$  by a DPM  $G(\theta)$ . Motivations for using a DPM prior include the possibility of accurately describing a large variety of prior information with this mixture model. The DPM classification model for class  $C_i$  can be written:

$$\theta_i \sim G, \quad (13)$$

$$\mathbf{x}_i \sim p(\mathbf{x}_i|\theta_i), \quad (14)$$

or, equivalently,

$$F \sim \mathcal{DP}(F_0, \alpha), \quad (15)$$

$$\phi_i \sim F, \quad (16)$$

$$\theta_i \sim f(\theta_i|\phi_i), \quad (17)$$

$$\mathbf{x}_i \sim p(\mathbf{x}_i|\theta_i), \quad (18)$$

where  $G(\theta)$  is a DPM over  $\Theta$ , with base distribution  $F_0$  and mixture distribution  $f$ . Learning a DPM classifier consists of estimating the DPM parameter  $\alpha$ , the parameters  $\theta_i$ , and the hyperparameters  $\phi_i$  from the training set  $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ .

**Example.** In order to illustrate these elements, we propose the following toy example. Assume that  $\theta$  is a random variable with Gaussian distribution, i.e.,  $f(\theta|\phi) = \mathcal{N}(\theta; \mu, \sigma^2)$ . The hyperparameter vector associated to this prior is  $\phi = [\mu, \sigma^2]$ . Assume also that the base distribution  $F_0$  is the product of a Gaussian distribution over  $\mu$  and an inverse gamma distribution over  $\sigma^2$ . Then, it can be shown that, conditional on the parameters of  $F_0$ ,  $\theta$  is distributed according to a mixture of Gaussians

$$p(\theta) = \sum_{k=1}^{\infty} \omega_k \mathcal{N}(\theta; \mu_k, \sigma_k^2), \quad (19)$$

where the weights are given by the stick-breaking rule and the parameters  $\phi_k = [\mu_k, \sigma_k^2]$  are distributed according to  $F_0$ , i.e., a Gaussian distribution for  $\mu$  and an inverse gamma distribution for  $\sigma^2$ . This defines the prior distribution over  $\theta$ . Note that the number of Gaussians defining the mixture is possibly infinite. In practice, only a finite number of Gaussians will have a significant weight, and the posterior distribution relies on these significant Gaussians. Therefore, the number of Gaussians is estimated in an implicit way. This example will be generalized to a multivariate parameter  $\theta$  in Section 5.

#### 4.1 Learning the Class Posterior Distribution

This section derives the posterior parameter distribution which follows from the model (15)-(18). By introducing the latent variables  $\mathbf{z}$  and cluster locations  $\mathbf{U}$ , one has

$$\begin{aligned} p(\theta|\mathbf{X}) &= \int \sum_{\mathbf{z}} p(\theta, \mathbf{z}, \mathbf{U}|\mathbf{X}) d\mathbf{U} \\ &= \int \sum_{\mathbf{z}, \mathbf{z}} f(\theta|U_z) p(\mathbf{z}|\mathbf{z}) p(\mathbf{U}, \mathbf{z}|\mathbf{X}) d\mathbf{U}, \end{aligned} \quad (20)$$

where  $\mathbf{z} = (z_1, \dots, z_N)$  is the indicator vector for  $(\theta_1, \dots, \theta_N)$ , whereas  $z$  is the indicator variable for the parameter  $\theta$  associated to the observation vector  $\mathbf{x}$  to be classified. The distribution  $p(\mathbf{z}|\mathbf{z})$  in (20) is defined as  $p(z = k|\mathbf{z}) = \frac{1}{N} \sum_{i=1}^N \delta_{k, z_i}$ , i.e., the number of latent variables  $z_i$  equal to  $k$  divided by the total number of latent variables (the number of samples in the training set). The parameter posterior distribution follows

$$\begin{aligned} p(\mathbf{U}, \mathbf{z}|\mathbf{X}) &= \int_{\Theta^N} p(\theta, \mathbf{U}, \mathbf{z}|\mathbf{X}) d\theta \\ &\propto \int_{\Theta^N} \mathcal{DP}(F_0, \alpha) \prod_{i=1}^N p(\mathbf{x}_i|\theta_i) f(\theta_i|U_{z_i}) d\theta_i, \end{aligned} \quad (21)$$

where  $\mathcal{DP}(F_0, \alpha)$  is represented by its latent variables and cluster locations  $(\mathbf{U}, \mathbf{z})$ . Merging Algorithms 1 and 2 yields the following sampling strategy referred to as Algorithm 3 below.

**Algorithm 3.** Sampling from the class posterior distribution

% Step 0: Initialization

- For  $i = 1, \dots, N$ , sample from the Polya urn  $\tilde{\phi}^{(0)} \sim p(\phi|\tilde{\phi}_1^{(0)}, \dots, \tilde{\phi}_{i-1}^{(0)})$  and deduce  $\tilde{\mathbf{z}}^{(0)}$  and  $\tilde{\mathbf{U}}^{(0)}$
- For  $i = 1, \dots, N$ , sample  $\tilde{\theta}_i^{(0)} \sim f(\theta|\tilde{\phi}_i^{(0)})$

% Step 1: Iterations For  $l = 1, \dots, L$ , do

- 1.1- Sample  $\tilde{\mathbf{z}}^{(l)} \sim p(\mathbf{z}|\tilde{\mathbf{U}}^{(l-1)}, \tilde{\theta}^{(l-1)})$  using Algorithm 4 or Algorithm 5

- 1.2- Sample  $\tilde{\mathbf{U}}^{(l)} \sim p(\mathbf{U}|\tilde{\mathbf{z}}^{(l)}, \tilde{\theta}^{(l-1)})$  using Algorithm 6  
 1.3- Sample  $\tilde{\theta}^{(l)} \sim p(\theta|\tilde{\mathbf{z}}^{(l)}, \tilde{\mathbf{U}}^{(l)})$ : for  $i = 1, \dots, N$ , sample

$$\tilde{\theta}_i^{(l)} \sim p(\theta_i|\mathbf{x}_i, \tilde{\mathbf{U}}_{z_i}^{(l)}) \propto p(\mathbf{x}_i|\theta_i)f(\theta_i|\tilde{\mathbf{U}}_{z_i}^{(l)}) \text{ using,} \\ \text{e.g., a MH step}$$

- 1.4- Sample  $\tilde{\xi}^{(l)} \sim p(\xi|\tilde{\mathbf{z}}^{(l)})$  such that  

$$p(\xi = k|\tilde{\mathbf{z}}^{(l)}) = \frac{1}{N} \sum_{i=1}^N \delta_{k, \tilde{z}_i^{(l)}}$$

- 1.5- Sample  $\tilde{\theta}^{(l)} \sim f(\theta|\tilde{\mathbf{U}}_{\xi^{(l)}}^{(l)})$

- 1.6- (optional) Sample the DPM hyperparameters, i.e.,  $\alpha$  and the parameters defining  $F_0$ .

The algorithms used to sample the latent variables and the cluster locations are detailed in Appendix B. An important extension is Step 1.6, which consists of defining a prior for the hyperparameter  $\alpha$  and the parameters of  $F_0$ . The parameter  $\alpha$  can be sampled using the method described in [21] which assumes a gamma prior for  $\alpha$ . The results of [21] show the robustness of the algorithm to  $\alpha$ . It is important to note that the convergence of the algorithm is made easier by initializing  $\alpha$  at a large value, and letting it decrease. This way, there are as many locations  $U_i$  as data  $\mathbf{x}_i$  in the beginning of the algorithm. Then, during the iterations, Step 1.6 updates the value of  $\alpha$ . At the very first iterations, this generally results in a decrease of its value (see Fig. 4), and therefore, the number of locations  $U_i$  also decreases. After some iterations, the value of  $\alpha$  stabilizes. This procedure ensures good convergence because we make sure that all possible modes of the final distributions are included at the beginning of the algorithm and that the less significant are removed along the iterations. This somehow comes down to following the regularization path, starting from the less regularized solution, toward the most regularized solutions, and stabilizing around the appropriate regularization level. The parameters of  $F_0$  can be sampled as shown in Section 5.

## 5 APPLICATION TO THE CLASSIFICATION OF SYNTHETIC AND REAL DATA

The efficiency of the proposed DPM classifier has been tested on different synthetic and real data sets. This section summarizes some experiments that illustrate the interest of the proposed classification strategy.

### 5.1 Classification of Radar Altimeter Signals

This section studies a problem, recently introduced in [24], consisting of classifying radar altimeter data. The objective is to validate the proposed DPM-based classification method on real data. The use of altimetry measurements on ocean surface and continental ice has demonstrated its effectiveness. However, the quality of the onboard algorithms also allows the acquisition of several measurements over non-ocean surfaces (coastal areas and inland water). The corresponding measurements are being considered bad, although they probably contain useful geophysical information. Consequently, new processing algorithms efficient for all kinds of surfaces need to be developed. A first identified step before developing such processing algorithms is the classification of altimeter waveforms, according to the over-flown surface. More precisely, this paper addresses the problem, recently introduced in [24], of classifying altimeter

waveforms in the four classes ‘‘Ocean,’’ ‘‘Desert,’’ ‘‘Forest,’’ and ‘‘Ice.’’ The proposed DPM-based classification algorithm will show very promising results for this problem.

#### 5.1.1 Statistical Model

Altimeter waveforms studied here represent the power of the backscattered altimeter echo. The observed time series are associated to the following statistical model:

$$\mathbf{x}(t) = M(t, \theta)\epsilon(t), \quad t = 1, \dots, T, \quad (22)$$

where  $\mathbf{x} = [\mathbf{x}(1), \dots, \mathbf{x}(T)]$  is the observed time series,  $T = 104$  is the number of samples, and  $\epsilon(t)$  is an i.i.d. multiplicative noise referred to as *speckle* in the radar community (see, for instance, [25, p. 96], [26] for more details). Note that the generative model introduced in (1) reduces here to

$$\mathbf{x}_{i,j}(t) = M_j(t, \theta_{i,j}, \epsilon_{i,j}) = M(t, \theta_{i,j})\epsilon_{i,j}(t), \quad (23)$$

where the model  $M_j = M$  is the same for all classes (the classes only differ by the distribution of the parameter vector  $\theta$ ) and where the dependence to the time instant  $t$  has been made explicit. The noise  $\epsilon(t)$  is known to be distributed according to a gamma distribution (see [27, p. 474] for the definition of the gamma distribution used in this paper), denoted as

$$\epsilon(t) \sim \text{Ga}(\epsilon(t); L, L), \quad (24)$$

where  $L$  is a known parameter (called *number of looks* in the radar community) associated to the altimeter. In this paper, we will use real data from the ENVISAT radar such that  $L = 100$ . The functional  $M(t, \theta)$  in (22) results from the so-called Brown model introduced in [28] and described accurately in [29, Section 2.7.1.1]:<sup>3</sup>

$$M(t, \theta) = P_n + \frac{a_\zeta \sigma_0}{2} \left[ 1 + \text{erf} \left( \frac{t - \tau - c_\zeta \sigma_c^2}{\sqrt{2} \sigma_c} \right) \right] \\ \exp \left[ -c_\zeta \left( t - \tau - \frac{c_\zeta \sigma_c^2}{2} \right) \right], \quad (25)$$

where

$$a_\zeta = \exp \left( \frac{-4 \sin^2 \zeta}{\xi_2} \right), \quad (26)$$

$$c_\zeta = \xi_1 \left[ \cos(2\zeta) - \frac{\sin^2(2\zeta)}{\xi_2} \right], \quad (27)$$

and where  $\tau$ ,  $\sigma_0$ ,  $\sigma_c$ , and  $\zeta$  are the unknown parameters,  $P_n$  is the known thermal noise level (it is estimated by using the first samples of the return waveform), and  $\xi_1, \xi_2$  are fixed coefficients (depending on the antenna bandwidth, the mean satellite altitude, the velocity of light, and the earth radius). The  $\text{erf}(\cdot)$  function is the error function defined by  $\text{erf}(u) = 2/\sqrt{\pi} \int_0^u \exp(-v^2) dv$ . The unknown parameters have the following physical interpretations (the interested reader is invited to consult [30] for more details):

3. The interested reader is invited to consult [29, Section 2.7.1.1] for more details. In particular, the full Hayne model includes additional terms which are negligible in the detection/classification problem, and are thus omitted here.

- $\tau > 0$  is the radar echo propagation time, related to the distance between the radar and the point observed on earth,
- $\sigma_0 > 0$  is the retrodiffusion coefficient depending on the ocean area,
- $\sigma_s > 0$  is related to the ocean significant wave height (SWH) by  $\text{SWH} = 2c_{\text{light}}\sigma_s$  with  $\sigma_s^2 = \sigma_c^2 - \sigma_p^2$  ( $c_{\text{light}}$  is the speed of light and  $\sigma_p^2$  is a known coefficient),
- $\zeta$  is the radar antenna absolute off-nadir pointing angle.

The previous statistical model is fully characterized by the following parameter vector:

$$\theta = [\tau/\tau^{\text{ref}}, \sigma_0/\sigma_0^{\text{ref}}, \text{SWH}/\text{SWH}^{\text{ref}}, \zeta/\zeta^{\text{ref}}] \in \Theta \subset \mathbb{R}^4,$$

where each parameter has been scaled to ensure the components of  $\theta$  are on a single scale (note that the dimension of the parameter space is  $d = 4$  for this example). The pdf of the radar waveform satisfies the following relation:

$$p(\mathbf{x}|\theta) = \prod_{t=1}^T \mathcal{G}a\left(\mathbf{x}(t); L, \frac{L}{M(t, \theta)}\right) \propto \exp\left[-L \sum_{t=1}^T \left[\log M(t, \theta) + \frac{\mathbf{x}(t)}{M(t, \theta)}\right]\right]. \quad (28)$$

### 5.1.2 Parameter and Hyperparameter Priors

The prior distribution for the unknown parameter vector  $\theta$  is a DPM with mixture distribution  $f(\theta|\phi)$ , base distribution  $F_0(\phi)$ , and DPM hyperparameter  $\alpha$ . The mixture distribution is assumed to be a multivariate Gaussian distribution,

$$f(\theta|\phi) = \mathcal{N}(\theta; \mu, \Sigma), \quad (29)$$

where  $\phi = \{\mu, \Sigma\}$  contains the hyperparameters ( $\mu$  is the mean vector and  $\Sigma$  is the diagonal covariance matrix). For computational reasons, we assume a normal-inverse Wishart base distribution as in [31],

$$F_0(\phi) = \mathcal{NIW}(\mu, \Sigma; \mu_0, \kappa_0, \nu_0, \Lambda_0), \quad (30)$$

$$= \mathcal{N}(\mu; \mu_0, \Sigma/\kappa_0) \mathcal{IW}(\Sigma; \nu_0, \Lambda_0), \quad (31)$$

where  $\mu_0$ ,  $\kappa_0$ ,  $\nu_0$ , and  $\Lambda_0$  are the base function hyperparameters (to which might be added the DPM hyperparameter  $\alpha$ ). The normal-inverse Wishart distribution is the conjugate prior for the mean and covariance matrix of a multivariate normal distribution. The reader is invited to consult [32, p. 272], [27, p. 80] for more details about this distribution. In particular, the inverse-Wishart distribution  $\mathcal{IW}(\Sigma; \nu_0, \Lambda_0)$  has the following density:

$$\mathcal{IW}(\Sigma; \nu_0, \Lambda_0) = \frac{2^{-\nu_0 d/2} \pi^{-d(d-1)/4}}{\prod_{i=1}^d \Gamma\left(\frac{\nu_0+1-i}{2}\right) |\Lambda_0|^{\nu_0/2} |\Sigma|^{(\nu_0+d+1)/2}} \exp\left[-\frac{1}{2} \text{trace}(\Lambda_0^{-1} \Sigma^{-1})\right]. \quad (32)$$

### 5.1.3 DPM Posterior Distribution

Before proposing the algorithm to sample the parameters and hyperparameters according to their full posterior, we first note that the Gaussian distribution is conjugate for the normal-inverse Wishart distribution. Indeed,

$$\begin{aligned} & \left[ \prod_{i=1}^N f(\theta_i|\phi) \right] F_0(\phi) \\ &= \left[ \prod_{i=1}^N \mathcal{N}(\theta_i; \mu, \Sigma) \right] \mathcal{NIW}(\mu, \Sigma; \mu_0, \kappa_0, \nu_0, \Lambda_0) \\ &= C(\theta_{1:N}) \mathcal{NIW}(\mu, \Sigma; \mu_N, \kappa_N, \nu_N, \Lambda_N), \end{aligned} \quad (33)$$

where

$$\bar{\theta} = \frac{1}{N} \sum_{i=1}^N \theta_i, \quad (34)$$

$$\mu_N = \frac{\kappa_0 \mu_0 + N \bar{\theta}}{\kappa_0 + N}, \quad (35)$$

$$\kappa_N = \kappa_0 + N, \quad (36)$$

$$\nu_N = \nu_0 + N, \quad (37)$$

$$\begin{aligned} \Lambda_N^{-1} &= \Lambda_0^{-1} + \sum_{i=1}^N (\theta_i - \bar{\theta})(\theta_i - \bar{\theta})^T \\ &+ \frac{\kappa_0 N}{\kappa_0 + N} (\mu_0 - \bar{\theta})(\mu_0 - \bar{\theta})^T, \end{aligned} \quad (38)$$

and the constant is (here,  $d = 4$ )

$$C(\theta_{1:N}) = \left[ \frac{\kappa_0}{\kappa_N} \right]^{d/2} \pi^{-\frac{Nd}{2}} \frac{|\Lambda_N|^{\nu_N/2}}{|\Lambda_0|^{\nu_0/2}} \prod_{j=1}^d \frac{\Gamma(\frac{\nu_N+1-j}{2})}{\Gamma(\frac{\nu_0+1-j}{2})}. \quad (39)$$

This constant depends on  $\theta_{1:N}$  through the parameter  $\Lambda_N$  and should be computed via its logarithm for numerical stability. Moreover, considering (33) for one  $\theta$ , we have

$$\int f(\theta|\phi) F_0(\phi) d\phi = C(\theta) \quad (40)$$

since the normal-inverse Wishart integrates to one.

### 5.1.4 Sampling the Hyperparameters

Algorithm 3 contains an optional step for estimating the hyperparameter  $\alpha$  of the DPM model and the hyperparameters defining the base distribution  $F_0(\phi)$ . This optional step requires to define appropriate prior distributions for these unknown hyperparameters. This paper assumes that the prior for  $\alpha$  is a gamma distribution as in [21]. The base distribution  $F_0(\phi)$  is a normal-inverse Wishart with parameters  $\mu_0$ ,  $\kappa_0$ ,  $\nu_0$ , and  $\Lambda_0$ . We assume here that  $\mu_0$  is known and we estimate  $\kappa_0$ ,  $\nu_0$ , and  $\Lambda_0$ . Therefore,

$$\begin{aligned} & p(\kappa_0, \nu_0, \Lambda_0 | \mathbf{U}, \mathbf{z}) \\ & \propto \left[ \prod_{z \in \mathcal{I}(z)} F_0(U_z | \kappa_0, \nu_0, \Lambda_0) \right] p(\kappa_0, \nu_0, \Lambda_0), \end{aligned} \quad (41)$$

where  $\mathcal{I}(z)$  denotes the set of values taken by the variables  $z_1, \dots, z_l$  and  $p(\kappa_0, \nu_0)$  is the prior distribution of the hyperparameters. Assume the prior structure

$$\begin{aligned} p(\kappa_0, \nu_0, \Lambda_0) &\propto \mathcal{G}a(\kappa_0; \alpha_\kappa, \beta_\kappa) \mathcal{P}(\nu_0; \lambda_\nu) \mathcal{IW}(\Lambda_0; \nu_{00}, \Lambda_{00}) \\ & \quad \prod_{\{d, d+1, \dots\}}(\nu_0), \end{aligned} \quad (42)$$



where  $\mathcal{P}(\nu_0; \lambda_\nu)$  is the Poisson distribution and  $\mathbb{I}_{\{d, d+1, \dots\}}(\cdot)$  is the indicator variable for the set  $(d, \infty)$ . The hyperparameter posterior distribution can be decomposed into the terms

$$p(\kappa_0 | \mathbf{U}, \mathbf{z}) = \mathcal{G}a(\kappa_0; \alpha'_\kappa, \beta'_\kappa), \quad (43)$$

$$p(\nu_0 | \Lambda_0, \mathbf{U}, \mathbf{z}) \propto \lambda'_\nu \Gamma(\nu_0 + 1)^{-1} \prod_{i=1}^d \Gamma\left(\frac{\nu_0 + 1 - i}{2}\right)^{-\#\mathcal{I}(\mathbf{z})} \mathbb{I}_{(d, \infty)}(\nu_0), \quad (44)$$

$$p(\Lambda_0 | \nu_0, \mathbf{U}, \mathbf{z}) = \mathcal{IW}(\Lambda_0; \nu', \Lambda'), \quad (45)$$

where  $\#\mathcal{I}(\mathbf{z})$  denotes the number of elements in  $\mathcal{I}(\mathbf{z})$  and

$$\alpha'_\kappa = \alpha_\kappa + \frac{d \#\mathcal{I}(\mathbf{z})}{2}, \quad (46)$$

$$\beta'_\kappa = \beta_\kappa + \frac{1}{2} \sum_{z \in \mathcal{I}(\mathbf{z})} (\mu_z - \mu_0)^T \Sigma_z^{-1} (\mu_z - \mu_0), \quad (47)$$

$$\lambda'_\nu = \lambda_\nu 2^{-\frac{d \#\mathcal{I}(\mathbf{z})}{2}} |\Lambda_0|^{-\frac{\#\mathcal{I}(\mathbf{z})}{2}} \prod_{z \in \mathcal{I}(\mathbf{z})} |\Sigma_z|^{-1/2}, \quad (48)$$

$$\nu' = \nu_0 \#\mathcal{I}(\mathbf{z}) + \nu_{00}, \quad (49)$$

$$\Lambda'^{-1} = \Lambda_{00}^{-1} + \sum_{z \in \mathcal{I}(\mathbf{z})} \Sigma_z^{-1}. \quad (50)$$

These hyperparameters can be sampled conditionally on  $\mathbf{U}$  and  $\mathbf{z}$  in Step 1.6 of Algorithm 3 by direct sampling for  $\kappa_0$  in (43) and  $\Lambda_0$  in (45), and by an MH step for  $\nu_0$  with truncated Poisson proposal distribution.

## 5.2 Experiments

This section first presents simulation results obtained with synthetic data, in order to demonstrate the accurate learning ability of the proposed Bayesian model. Results obtained with real radar altimetry data are then analyzed.

### 5.2.1 Synthetic Data

The data studied here have been generated using the Brown model in (25). In order to demonstrate the good ability of our model to learn the posterior data distribution for a class, we consider one class with the following distribution for the model parameters (adjusted to match real data):

$$\tau \sim \mathcal{N}(\tau; 32, 4 \cdot 10^{-4}) \text{ and } \tau^{\text{ref}} = 32, \quad (51)$$

$$\sigma_0 \sim \mathcal{N}(\sigma_0; 12, 2.25) \text{ and } \sigma_0^{\text{ref}} = 12, \quad (52)$$

$$\text{SWH} \sim 0.5 \mathcal{G}a(\text{SWH} - 4; 10, 10) + 0.5 \mathcal{U}(\text{SWH}; [1, 2]) \text{ and } \text{SWH}^{\text{ref}} = 2, \quad (53)$$

$$\zeta \sim \mathcal{N}(\zeta; 9.5 \cdot 10^{-4}, 6.3 \cdot 10^{-9}) \text{ and } \zeta^{\text{ref}} = 0.001. \quad (54)$$

Fig. 1 shows 10 waveforms generated according to the Brown model in (25), with parameters sampled from the distributions (51)-(54) with multiplicative gamma-distributed white noise ( $L = 100$ ). Step 1.3 of the above algorithm

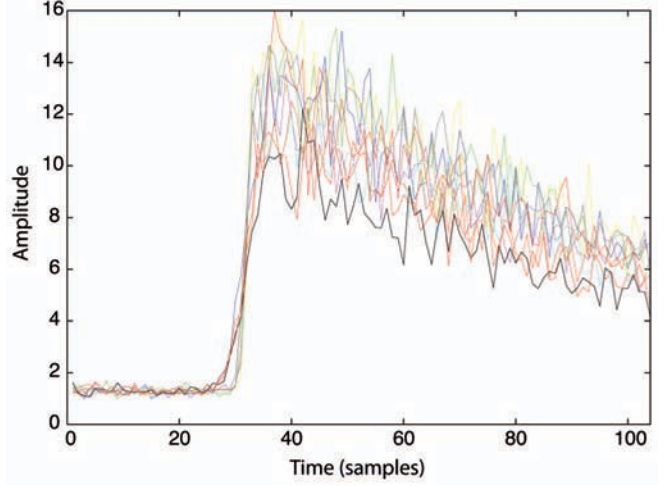


Fig. 1. Ten of the simulated radar altimetry time series (i.e., 10 data) used to train the classifier.

has been implemented using a one-at-a-time Gaussian random walk with variance  $4 \cdot 10^{-4}$ . The model hyperparameters have been selected as follows:  $\mu_0 = [1 \ 1 \ 1]^T$ ,  $\alpha_\kappa = \beta_\kappa = 5$ ,  $\lambda_\nu = d = 4$ ,  $\Lambda_{00} = 0.1 \mathbf{I}_d$ , and  $\nu_{00} = d + 1$ . These values are generic and changing them into a reasonable range does not much influence the results.

In order to check that the algorithm accurately learns the values of the different parameters for each training data set, the proposed algorithm has been run 100 times for  $N = 100$  different training data sets. Each Monte Carlo run consists of  $L = 5,000$  iterations (including  $L_{\text{burn in}} = 3,000$  iterations). The minimum mean square error (MMSE) estimates of  $\theta_i^{\text{gt}}$ , denoted as  $\hat{\theta}_i^{\text{MMSE}}$ , have been compared to their ground truth values via

$$\text{Err}_\theta(N) = \frac{1}{N} \sum_{i=1}^N (\hat{\theta}_i^{\text{MMSE}} - \theta_i^{\text{gt}})^T \mathbf{S}^{-1} (\hat{\theta}_i^{\text{MMSE}} - \theta_i^{\text{gt}}), \quad (55)$$

where  $\theta_i^{\text{gt}}$  is the actual value of the parameter vector,  $\mathbf{S}^{-1} = \text{diag}(\tau^{\text{ref}}, \sigma_0^{\text{ref}}, \text{SWH}^{\text{ref}}, \zeta^{\text{ref}})$ , and

$$\hat{\theta}_i^{\text{MMSE}} = \frac{1}{L - L_{\text{burn in}}} \sum_{l=L_{\text{burn in}}+1}^L \tilde{\theta}_i^{(l)}, \quad (56)$$

where  $\tilde{\theta}_1 = \tau/\tau^{\text{ref}}$ ,  $\tilde{\theta}_2 = \sigma_0/\sigma_0^{\text{ref}}$ ,  $\tilde{\theta}_3 = \text{SWH}/\text{SWH}^{\text{ref}}$ , and  $\tilde{\theta}_4 = \zeta/\zeta^{\text{ref}}$  are the normalized model parameters. The average of  $\text{Err}_\theta(N)$  over the 100 Monte Carlo runs is 0.109 with standard deviation 0.063 (to be compared with the average norm of  $\theta_i^{\text{gt}}$  over the 100 simulations, which is 5.43). Fig. 2 shows that the parameters of each training signal have been accurately learned through the simulations. Over the 100 simulations, only six failed to converge to the ground truth parameter values for each of the 100 parameters. In order to illustrate the algorithm convergence, Figs. 3 and 4 represent the simulated components of  $\theta$  sampled in Step 1.5, the evolution of  $\#\mathcal{I}(\mathbf{z})$ , and the hyperparameters  $\alpha$ ,  $\nu$ , and  $\kappa$  for one typical simulation with  $N = 100$  data. The chain has clearly converged after about 3,000 burn-in iterations. Overall, the learning phase for this example takes a couple of minutes on a standard PC computer and Matlab/Octave code. Fig. 5 represents the histograms of the four components of the samples  $\tilde{\theta}^{(l)}$  for  $l \in \{3,001, \dots, 5,000\}$ . The

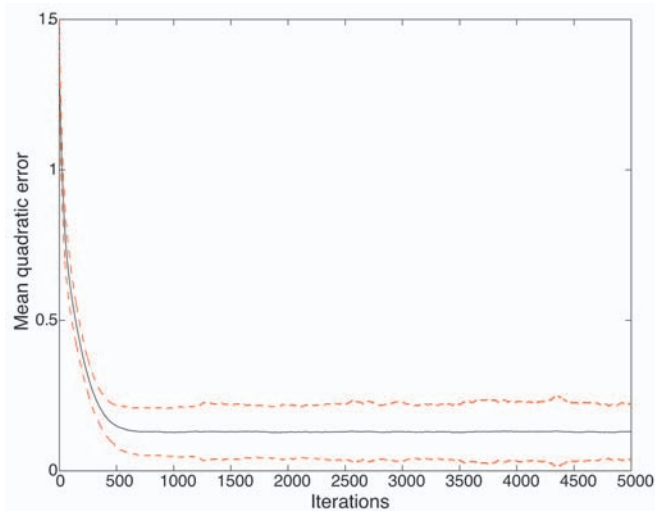


Fig. 2. Evolution of the mean quadratic error for the signal parameters  $\theta^{(l)}$ ,  $l = 1, \dots, 100$  for each of the 5,000 iterations of the MCMC algorithm, averaged over 100 Monte Carlo runs. The upper and lower curves represent the 95 percent confidence interval ( $\pm 2$  standard deviations).

estimated marginal posterior distributions are quite close to the true parameter distributions, which shows the ability of our model to learn multimodal, non-Gaussian distributions.

Fig. 6 shows the marginal posterior distributions of the different parameters for two sample sizes,  $N = 10$  and  $N = 1,000$ , which have to be compared with Fig. 5, which has been obtained with  $N = 100$ . The higher  $N$ , the better the estimated parameter posteriors, as expected. When a very limited number of training data is available ( $N = 10$ ), the learning ability is also correct. The regularization role of the DPM during learning can be clearly seen on the third component of  $\theta$ : The model considers that there is only one mode for  $N = 10$  (which is reasonable given the very limited number of training data) whereas the two modes are clearly recovered for  $N = 100$  and  $N = 1,000$ . These simulations have clearly demonstrated the ability of the proposed algorithm to learn the parameter posterior distribution. The algorithm classification ability is then studied by using real radar altimetry data.

### 5.2.2 Real Data

As explained in [24], extraction of geophysical information from waveforms backscattered from nonoceanic surfaces is an identified goal of future altimetry missions. The classification of the received radar waveforms into the different classes

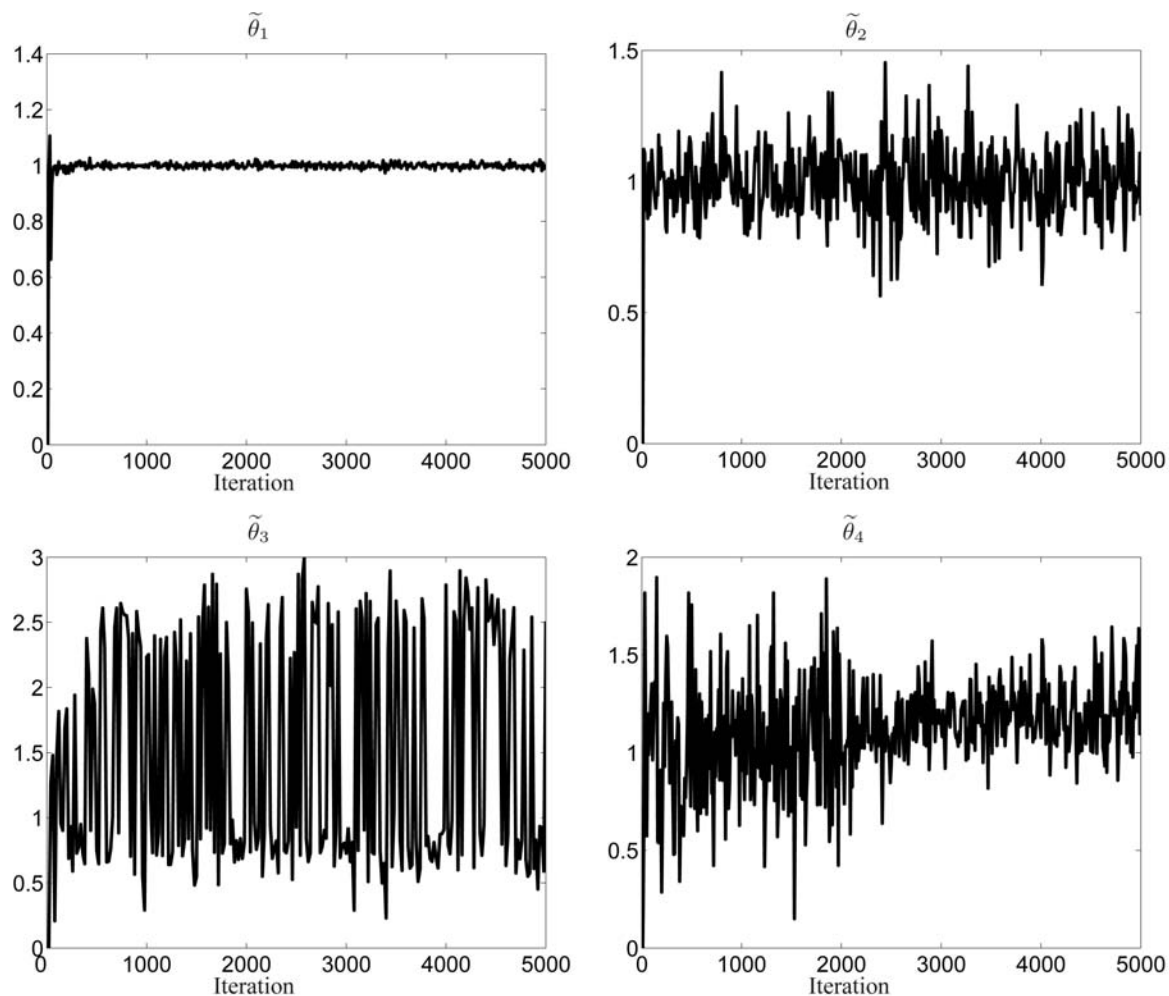


Fig. 3. Evolution of the parameter  $\theta$  sampled at each iteration in Step 1.5 of Algorithm 3 from the class posterior distribution, for each of the four components. As can be seen, the MCMC algorithm mixes quite well.

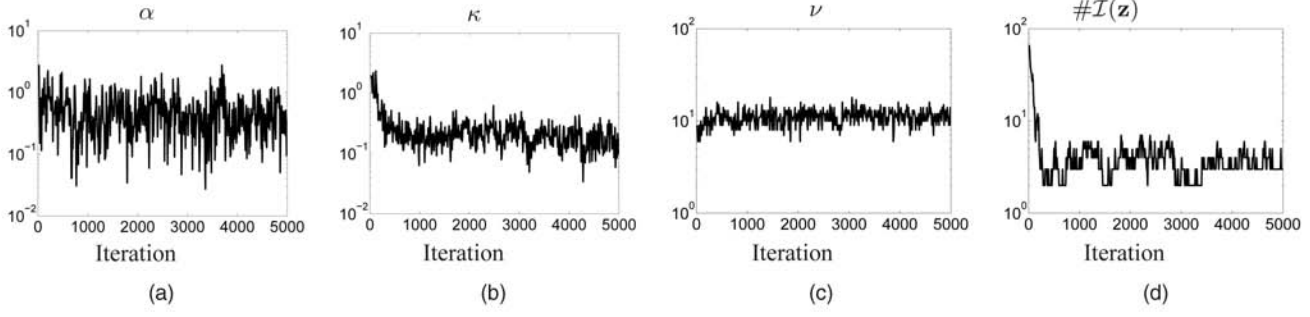


Fig. 4. Evolution of the hyperparameters  $\alpha$  ((a) in log scale),  $\kappa$  ((b) in log scale),  $\nu$  ((c) in log scale) sampled at each iteration in Step 1.6 of Algorithm 3. (d) represents the evolution of the number of different values for  $z$  used at each iteration (in log scale).

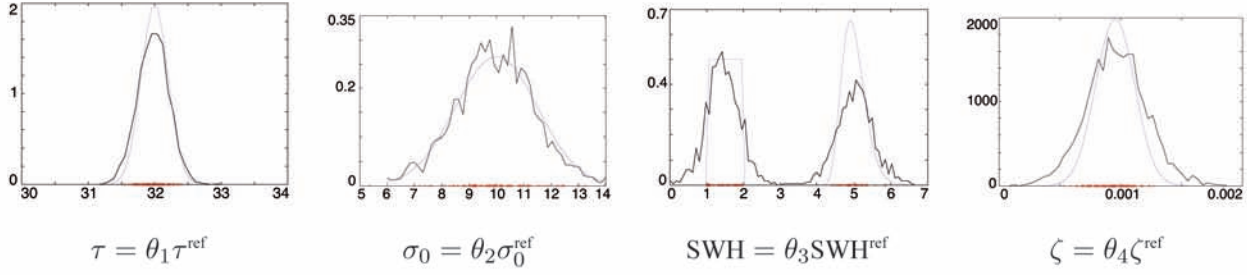


Fig. 5. Histograms of the four components of the parameter vector computed from the samples  $\tilde{\theta}^{(l)}$  for  $l \in \{3,001, \dots, 5,000\}$  (jagged lines), and pdfs used to sample the four components of the parameter vector (smooth lines). The crosses represent the values used in the simulations.

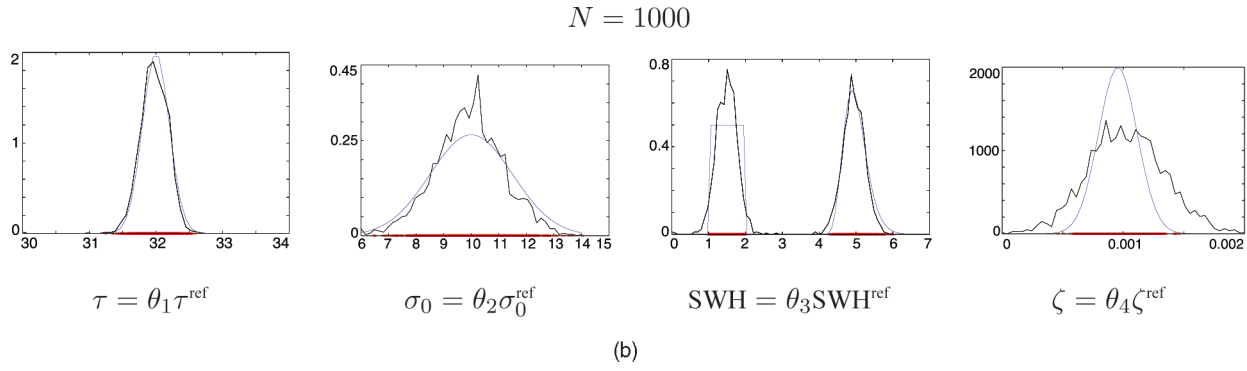
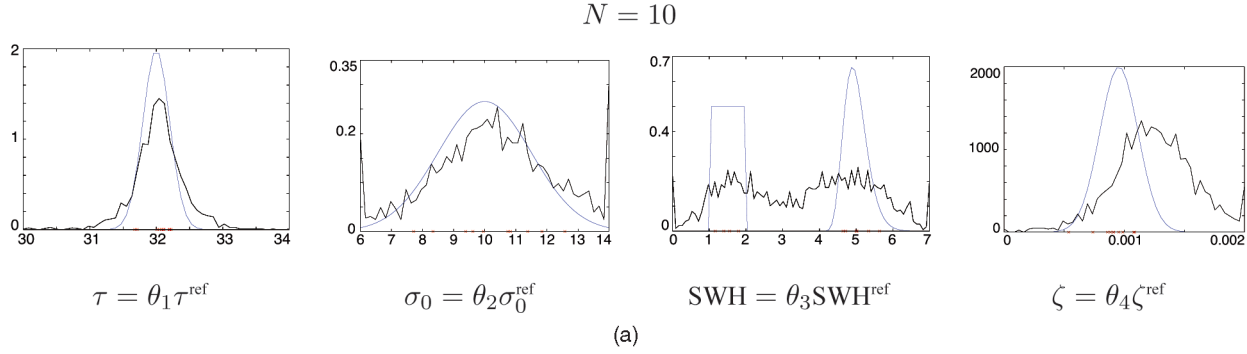


Fig. 6. Histograms of the four components of the parameter vector computed from the samples  $\tilde{\theta}^{(l)}$  for  $l \in \{3,001, \dots, 5,000\}$  (jagged lines), and pdfs used to sample the four components of the parameter vector (smooth lines). The crosses indicate the values used in the simulations. The number of training data is (a)  $N = 10$  and (b)  $N = 1,000$  (bottom row).

corresponding to the different kinds of the surfaces would considerably simplify this extraction of information. The real data studied in this section are grouped into four classes, referred to as “Ocean,” “Ice,” “Forest,” and “Desert” as described in Table 1. Ten training data from each class are plotted in Fig. 7, where it should be noted that the amplitudes

can differ significantly from one class to another, some having a large variability.

In order to assess the learning ability of our classification strategy,  $N = 1,000$  training data from each of the four classes were selected randomly for training. The MCMC algorithm was run independently on these four data sets, with  $\tau^{\text{ref}} = 32$ ,  $\sigma_0^{\text{ref}} = 9$ ,  $\text{SWH}^{\text{ref}} = 2$ , and  $\zeta^{\text{ref}} = 0.001$ . This

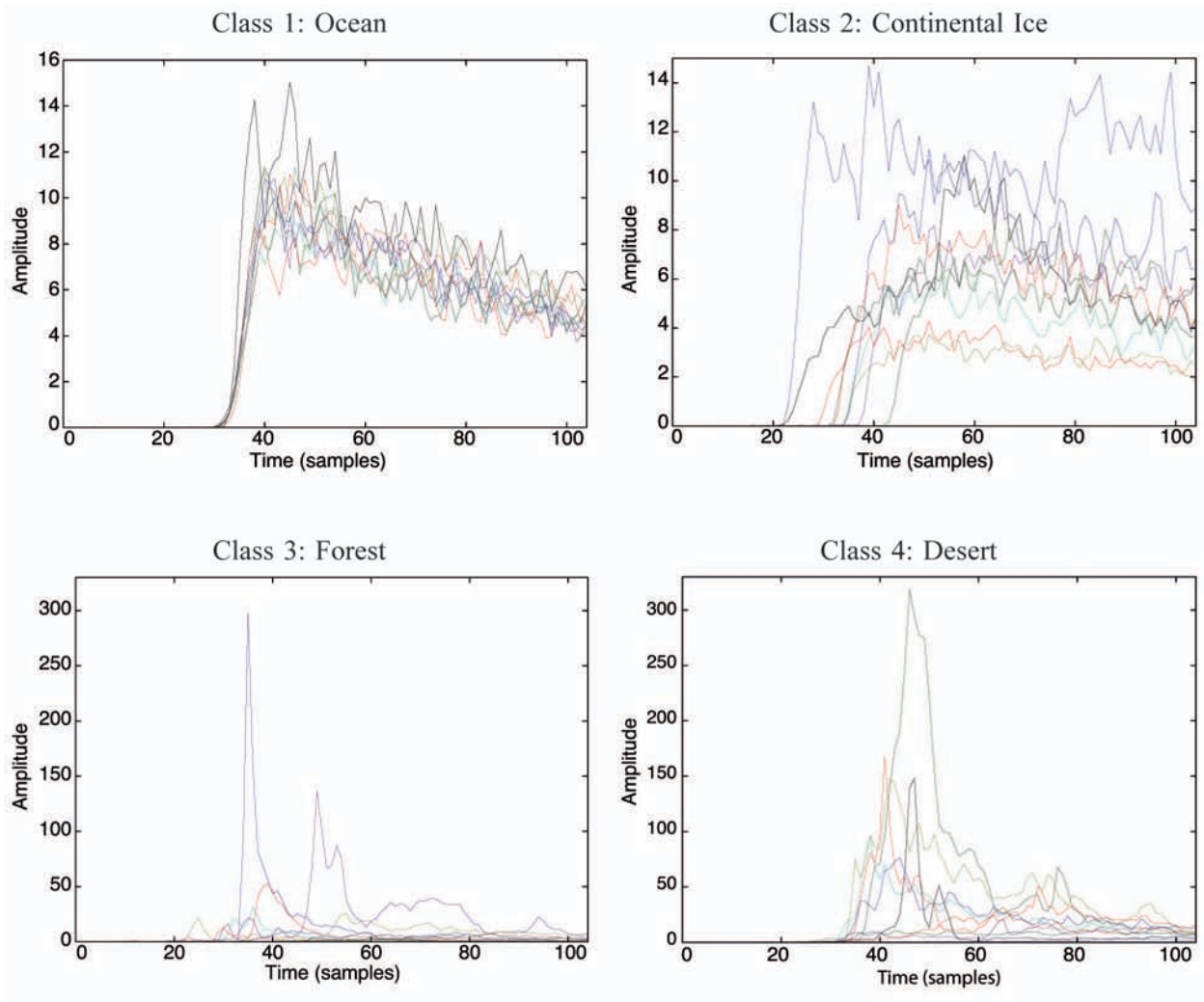


Fig. 7. Ten of the real altimetry time series (i.e., 10 data) used to train the classifier observed over the different surfaces.

procedure was repeated five times by changing the random sample of  $N = 1,000$  data for each class, leading to the histograms in Fig. 8. As can be seen, the procedure ensures a stable learning of the parameter posterior distribution. Moreover, it can be seen that this distribution is multimodal and/or non-Gaussian in several parameters, which justifies the approach implemented here. Once the model parameters have been estimated conditioned on each class, the estimates can be plugged into the class pdfs following the *plug-in* classification rule [4]. The results obtained with this plug-in rule for the four-class problem defined in Table 1 are depicted in Table 2.

For comparison, the results obtained when using a Gaussian prior for the unknown parameter vector are provided in Table 3. These results have been computed with the same code as in the DPM prior case with the

additional constraint  $\alpha = 10^{-300}$ . The proposed classification rule based on a DPM prior for  $\theta$  yields better performance than a naive Bayesian rule based on a Gaussian prior for  $\theta$ , although the latter is already quite sophisticated. The proposed DPM-based classifier was also compared to the one-versus-one and one-versus-all SVM multiclass methods with a polynomial kernel [33, p. 59]. The average of 50 confusion matrices obtained with 1,000 training samples and 1,000 test samples is provided in Tables 4 and 5. The DPM-based classification method shows better performance than the SVM approaches (because it uses the generative Brown model) at the price of a higher computational cost. The advantage of using a generative model allows one to encode the mapping between a reduced set of parameters and the data. An alternative strategy could be to estimate the parameters and thus implement an SVM-based classifier in the space of the parameters. However, this strategy is less robust insofar as it does not capture the uncertainty about the parameters.

TABLE 1  
Data Used to Train and Test the Classifier

Class	1: Ocean	2: Ice	3: Forest	4: Desert
Nb. of radar waveforms	9640	10480	8920	11000

## 6 CONCLUSIONS

This paper studied a new supervised classification algorithm based on a Bayesian model whose parameter priors

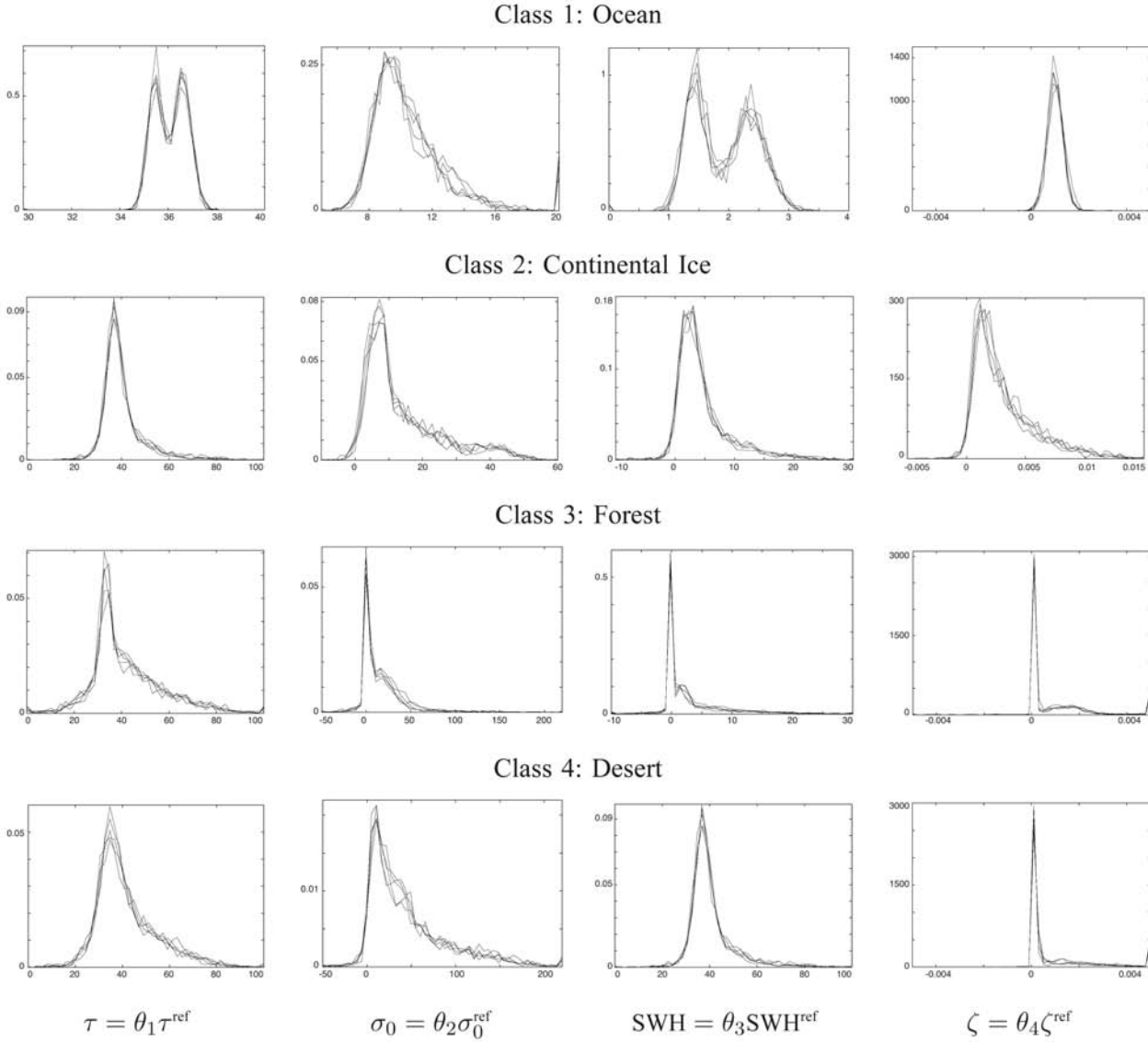


Fig. 8. Histograms of the four components of  $\theta$  computed from the samples  $\tilde{\theta}^{(l)}$  for  $l \in \{3,001, \dots, 5,000\}$  for the four classes (marginal distributions). In order to assess the stability of the learning procedure, five different training sets were used, each containing 1,000 training data.

are Dirichlet process mixtures. Note again that Dirichlet process mixtures are interesting since they can be used to describe accurately a large class of probability distributions. An efficient learning algorithm was developed to estimate the class posterior distributions required for Bayesian classification. The resulting classification rule was applied to a practical problem consisting of classifying radar altimeter waveforms backscattered from different surfaces

(i.e., oceans, ices, forests, and deserts). The obtained results are promising for this application. A recent study has shown that linear discriminant analysis performed on appropriate features (extracted from the Brown model) can also provide interesting classification results for altimetric signals [24]. A comparison between this classification strategy and the proposed DPM-based classification rule is currently under investigation.

TABLE 2  
Confusion Matrix for the Four-Class Altimetry Problem Using a DPM Prior for  $\theta$

	Estimated Class: 1	Estimated Class: 2	Estimated Class: 3	Estimated Class: 4
True class: 1	98.71%	1.08%	0.17%	0.04%
True class: 2	1.89%	65.79%	18.76%	13.56%
True class: 3	0.44%	12.47%	65.33%	21.76%
True class: 4	0.81%	18.96%	25.05%	55.17%

TABLE 3  
Confusion Matrix for the Four-Class Altimetry Problem Using a Gaussian Prior for  $\theta$

	Estimated Class: 1	Estimated Class: 2	Estimated Class: 3	Estimated Class: 4
True class: 1	97.08%	1.55%	0.17%	1.20%
True class: 2	8.13%	58.48%	20.03%	13.36%
True class: 3	4.04%	39.69%	35.19%	21.08%
True class: 4	1.93%	29.49%	20.68%	47.89%

TABLE 4  
Confusion Matrix for the Four-Class Altimetry Problem Using the One-versus-One SVM Multiclass Method [33, p. 59]  
(1,000 Training Samples and 1,000 Test Samples)

	Estimated Class: 1	Estimated Class: 2	Estimated Class: 3	Estimated Class: 4
True class: 1	95.91%	2.68%	0.38%	1.03%
True class: 2	34.18%	44.06%	11.12%	10.64%
True class: 3	4.82%	25.79%	44.52%	24.87%
True class: 4	8.23%	27.96%	21.24%	42.57%

TABLE 5  
Confusion Matrix for the Four-Class Altimetry Problem Using the One-versus-All SVM Multiclass Method [33, p. 59]  
(1,000 Training Samples and 1,000 Test Samples)

	Estimated Class: 1	Estimated Class: 2	Estimated Class: 3	Estimated Class: 4
True class: 1	59.76%	31.66%	5.95%	2.63%
True class: 2	18.85%	49.98%	19.52%	11.65%
True class: 3	3.67%	27.74%	47.64%	20.95%
True class: 4	4.34%	26.84%	23.71%	45.11%

## APPENDIX A

### A GLOSSARY OF NOTATIONS

- $J$ : number of classes
- $C_1, \dots, C_J$ : classes
- $\mathbf{X}_j = \{\mathbf{x}_{1,j}, \dots, \mathbf{x}_{N_j,j}\}$ : training set for class  $C_j$
- $\mathbf{x}_{i,j}$ :  $i$ th training vector for class  $C_j$
- $N_j$ : number of training vectors for class  $C_j$
- $M_j$ : model for class  $C_j$
- $\theta_{i,j}$ : parameter vector associated to the training vector  $\mathbf{x}_{i,j}$
- $\epsilon_{i,j}$ : noise associated to the training vector  $\mathbf{x}_{i,j}$
- $\Theta_j \subset \mathbb{R}^d$ : parameter space for class  $C_j$
- $d$ : dimension of the parameter space
- $p_j(\mathbf{x}_{i,j}|\theta_{i,j})$ : pdf of  $\mathbf{x}_{i,j}$  (belonging to the class  $C_j$ ) conditioned upon  $\theta_{i,j}$
- $p(\theta_j|\mathbf{X}_j)$ : parameter posterior distribution for the training set  $\mathbf{X}_j$  associated to class  $C_j$
- $p_j(\theta_j|\phi_j)$ : prior distribution for  $\theta_j$
- $\phi_j$ : hyperparameter vector for class  $C_j$
- $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_N)$ : generic training set
- $\mathbf{x}$ : generic observation vector to be classified
- $p(\mathbf{x}|\mathbf{X}_j)$ : predictive density of  $\mathbf{x}$  conditioned to the training set  $\mathbf{X}_j$
- $N$ : number of training vectors belonging to  $\mathbf{X}$
- $\theta$ : parameter vector for the generic training set  $\mathbf{X}$
- $p(\theta|\mathbf{X})$ : parameter posterior distribution for the generic training set  $\mathbf{X}$
- $\tilde{\theta}^{(1)}, \dots, \tilde{\theta}^{(L)}$ : samples distributed according to  $p(\theta|\mathbf{X})$
- $\phi$ : hyperparameter vector for the generic training set  $\mathbf{X}$
- $p(\phi|\mathbf{X})$ : hyperparameter posterior distribution for the generic training set  $\mathbf{X}$
- $\tilde{\phi}^{(1)}, \dots, \tilde{\phi}^{(L)}$ : samples distributed according to  $p(\phi|\mathbf{X})$
- $F \sim \mathcal{DP}(\alpha, F_0)$ : Dirichlet process with base distribution  $F_0$  and concentration parameter  $\alpha$
- $\omega_1, \omega_2, \dots$ : weights of the Dirichlet process (DP)
- $U_1, U_2, \dots$ : cluster locations of the Dirichlet process
- $\mathcal{B}(\omega; a, b)$ : beta distribution with parameters  $a$  and  $b$  for the random variable  $\omega$
- $f(\theta|\phi)$ : mixture density of a Dirichlet process mixture (DPM)
- $\mathbf{z} = (z_1, \dots, z_N)$ : vector containing the latent variables associated to the vectors  $\theta_1, \dots, \theta_N$  where  $\theta_i \sim G$
- $\mathbf{U} = (U_1, \dots, U_N)$ : cluster locations of  $\theta_1, \dots, \theta_N$  where  $\theta_i \sim G$
- $z$ : latent variable associated to a generic parameter vector  $\theta$

- $p(z|F)$ : discrete distribution of the latent variable  $z$  associated to the Dirichlet process  $F$
- $G(\theta)$ : DPM probability measure
- $M(t, \theta)$ : Model for an altimetric signal with parameter vector  $\theta$
- $\mathbf{x} = [\mathbf{x}(1), \dots, \mathbf{x}(T)]$ : time-varying radar waveform
- $\epsilon(t)$ : vector containing noise samples for an altimetric signal
- $\mathcal{G}a(\omega; a, b)$ : gamma distribution with parameters  $a$  and  $b$  for the random variable  $\omega$
- $P_n$ : known thermal noise level
- $\xi_1, \xi_2$ : known coefficients depending on the kind of altimeter radar
- $\tau, \text{SWH}, \sigma_e, \zeta$ : unknown parameters for the altimetric waveform
- $\mathcal{N}(\theta; \mu, \Sigma)$ : multivariate normal distribution with mean vector  $\mu$  and covariance matrix  $\Sigma$  for the random vector  $\theta$
- $\mathcal{NIW}(\mu, \Sigma; \mu_0, \kappa_0, \nu_0, \Lambda_0)$ : normal-inverse Wishart distribution with parameters  $\mu_0, \kappa_0, \nu_0, \Lambda_0$  for the random vector and matrix  $\mu$  and  $\Sigma$
- $\mathcal{IW}(\Sigma; \nu_0, \Lambda_0)$ : inverse Wishart distribution with parameters  $\nu_0, \Lambda_0$  for the random matrix  $\Sigma$
- $\mathcal{P}(\nu_0; \lambda_\nu)$ : Poisson distribution with parameter  $\lambda_\nu$  for the random variable  $\nu_0$
- $\mathcal{I}(\mathbf{z})$ : set of values taken by the variables contained in  $\mathbf{z}$
- $\#\mathcal{I}(\mathbf{z})$ : number of elements in  $\mathcal{I}(\mathbf{z})$
- $\mu, \Sigma$ : hyperparameters associated to  $\theta$  for the classification of altimetric signals
- $\mu_0, \kappa_0, \nu_0, \Lambda_0$ : base function hyperparameters associated to the hyperparameter prior
- $\mu_N, \kappa_N, \nu_N, \Lambda_N$ : parameters of the DPM posterior distribution for the classification of altimetric signals
- $\mathbb{I}_{\{d, d+1, \dots\}}(\cdot)$ : indicator variable for the set  $(d, \infty)$
- $\alpha'_\kappa, \beta'_\kappa, \lambda'_\nu, \nu', \Lambda'^{-1}$ : parameters of the hyperparameter posterior distribution for the classification of altimetric signals

## APPENDIX B

### SAMPLING DPM LATENT VARIABLES AND CLUSTER LOCATIONS

This appendix summarizes three algorithms which can be used to sample DPM latent variables and cluster locations:

- The Algorithm 4 details Step 1.1 of the Gibbs sampler in Algorithm 3 used to sample the DP in the case where the base distribution  $F_0$  is conjugate for the mixture distribution  $f(\cdot)$ , following [34, Algorithm 2].
- Algorithm 5 extends Algorithm 4 to the more general case where  $F_0$  is not conjugate for  $f(\cdot)$  and is derived from [34, Algorithm 5].
- Algorithm 6 addresses Step 1.2 of Algorithm 3 above in all cases.

Before presenting these algorithms, let us introduce some notations.

Let  $\mathcal{I}(\mathbf{z})$  denote the set of values taken by the variables  $z_1, \dots, z_N$ . Indeed, for the sake of simplicity, we do not assume that the  $z_i$ s take all their values in a set  $\{1, \dots, k_{\max}\}$ , but instead some values spread among the integers. The

reason for this is that the Gibbs samplers may create and suppress locations  $U_k$ s, thus incrementing the largest  $k$  while forming gaps in-between “alive” indexes  $k$ . Thus,  $\mathcal{I}(\mathbf{z})$  contains the values of  $k$  that are effectively used at a given iteration. We denote  $N_{-i,k}(\mathbf{z}) = \sum_{i'=1, i' \neq i}^N \delta_{k, z_{i'}}$  the number of variables  $z_{i'}$ s ( $i' \neq i$ ) which are equal to  $k$ .

**Algorithm 4.** Sampling the latent variables (conjugate case)

- Let  $\mathbf{z}' \leftarrow \tilde{\mathbf{z}}^{(l-1)}$
- For  $i = 1, \dots, N$ , do
  - if  $z'_i \neq z'_{i'}$  for all  $i' \neq i$  (singleton cluster), then remove  $\tilde{U}_i^{(l)}$  from the state.
  - Sample  $z'_i \sim p(z'_i | \mathbf{z}'_{-i}, \tilde{\theta}_i^{(l-1)}, \tilde{\mathbf{U}}^{(l-1)})$  such that

$$p(z'_i = k | \mathbf{z}'_{-i}, \tilde{\theta}_i^{(l-1)}, \tilde{\mathbf{U}}^{(l-1)}) \propto \begin{cases} N_{-i,k} f(\tilde{\theta}_i^{(l-1)} | \tilde{U}_k^{(l-1)}), & \text{for } k \in \mathcal{I}(\mathbf{z}'), \\ \alpha \int f(\tilde{\theta}_i^{(l-1)} | U) F_0(dU), & \text{for a new } k \notin \mathcal{I}(\mathbf{z}'), \\ 0, & \text{for all other values of } k, \end{cases} \quad (57)$$

and, if a new cluster is created, sample  $\tilde{U}_{z'_i}^{(l-1)} \sim p(U | \tilde{\theta}_i^{(l-1)}) \propto f(\tilde{\theta}_i^{(l-1)} | U) F_0(U)$

- Let  $\tilde{z}_i^{(l)} = z'_i$

There are many cases where the base distribution is not conjugate for  $f(\cdot)$  in practical applications. In this case, MH steps are necessary, and they may be implemented as follows (derived from [34, Algorithm 5]):

**Algorithm 5.** Sampling the latent variables (general case)

- Let  $\mathbf{z}' \leftarrow \tilde{\mathbf{z}}^{(l-1)}$
- For  $i = 1, \dots, N$ , do
  - Sample a candidate  $z_i^*$  from the Polya urn probabilities
$$Q(z_i^* = k) = \begin{cases} \frac{N_{-i,k}(\mathbf{z}')}{\alpha + N - 1}, & \text{for } k \in \mathcal{I}(\mathbf{z}'), \\ \frac{\alpha}{\alpha + N - 1}, & \text{for a new } k \notin \mathcal{I}(\mathbf{z}'), \\ 0, & \text{for all other values of } k. \end{cases} \quad (58)$$
  - if  $z_i^* \in \mathcal{I}(\mathbf{z}')$ , then compute  $r_1 = f(\tilde{\theta}_i^{(l-1)} | U_{z_i^*}') / f(\tilde{\theta}_i^{(l-1)} | U_{z_i}^{(l-1)})$ . With probability  $\min(1, r_1)$ , set  $z_i^{(l)} = z_i^*$ .
  - Otherwise, compute  $r_1 = f(\tilde{\theta}_i^{(l-1)} | U^*) / f(\tilde{\theta}_i^{(l-1)} | U_{z_i}^{(l-1)})$ , where  $U^* \sim F_0$ . With probability  $\min(1, r_1)$ , set  $z_i^{(l)} = z_i^*$  and  $U_{z_i}^{(l-1)} \leftarrow U^*$ .
  - Let  $\tilde{z}_i^{(n)} = z_i^{(l)}$

Once the latent variables are sampled, the cluster locations can be updated as indicated below.

**Algorithm 6.** Sampling the cluster locations (all cases)

- Let  $\mathbf{U}' \leftarrow \tilde{\mathbf{U}}^{(l-1)}$
- For  $k \in \mathcal{I}(\tilde{\mathbf{z}}^{(l)})$ , sample  $\tilde{U}_k^{(l)}$  from  $p(U_k | \tilde{\theta}_i^{(l-1)})$  with  $i$  such

$$\text{that } \tilde{z}_i^{(l)} = k \propto F_0(U_k) \prod_{\substack{i=1, \dots, N \\ \text{such that } \tilde{z}_i^{(l)} = k}} f(\tilde{\theta}_i^{(l-1)} | U_{-k}),$$

as follows:

- In the conjugate case, sample directly from  $p(U_k | \tilde{\theta}_i^{(l-1)})$  with  $i$  such that  $\tilde{z}_i^{(l)} = k$
- In the non-conjugate case, sample  $U^* \sim q(U^* | U_k')$ , then compute

$$r_2 = \frac{q(U'_k|U^*) F_0(U^*)}{q(U^*|U'_k) F_0(U'_k)} \prod_{\substack{i=1, \dots, N \\ \text{such that } \tilde{z}_i^{(l)} = k}} \frac{f(\tilde{\theta}_i^{(l-1)}|U^*)}{f(\tilde{\theta}_i^{(l-1)}|U'_k)}, \quad (59)$$

and, with probability  $\min(1, r_2)$ , set  $\tilde{U}_i^{(l)} \leftarrow U^*$ ;  
otherwise, set  $\tilde{U}_i^{(l)} \leftarrow U'_k$

The implementation of Algorithm 6 requires sampling from the normal-inverse Wishart distribution  $\mathcal{N}\mathcal{I}\mathcal{W}(\mu, \Sigma; \tilde{\mu}_N, \tilde{\kappa}_N, \tilde{\nu}_N, \tilde{\Lambda}_N)$ , where  $\tilde{\mu}_N$ ,  $\tilde{\nu}_N$ ,  $\tilde{\kappa}_N$ , and  $\tilde{\Lambda}_N$  are defined as in (34)-(38), with the  $\theta_i$ s being the  $\tilde{\theta}_i^{(l-1)}$ s such that  $\tilde{z}_i^{(l)} = k$ .

## ACKNOWLEDGMENTS

The authors would like to thank L. Amarouche, P. Thibaut, and O.Z. Zanife from *Collecte Localisation Satellites* (CLS), Toulouse, France, for providing the real data sets and for fruitful discussions regarding altimetry.

## REFERENCES

- [1] L. Devroye, L. Györfi, and G. Lugosi, *A Probabilistic Theory of Pattern Recognition*. Springer-Verlag, 1996.
- [2] R.O. Duda, P.E. Hart, and D.G. Stork, *Pattern Classification*, second ed. Wiley, 2001.
- [3] K. Fukunaga, *Introduction to Statistical Pattern Recognition*. Academic Press, 1990.
- [4] A.K. Jain, R.P.W. Duin, and J. Mao, "Statistical Pattern Recognition: A Review," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 22, no. 1, pp. 4-37, Jan. 2000.
- [5] S. Theodoridis and K. Koutroubas, *Pattern Recognition*. Academic Press, 1999.
- [6] R. Kass and A. Raftery, "Bayes Factors," *J. Am. Statistical Assoc.*, vol. 90, no. 430, pp. 773-795, 1995.
- [7] R. Raina, Y. Shen, A. Ng, and A. McCallum, "Classification with Hybrid Generative/Discriminative Models," *Advances in Neural Information Processing Systems*, vol. 16, MIT Press, 2004.
- [8] A. McCallum and K. Nigam, "A Comparison of Event Models for Naive Bayes Text Classification," *Proc. Nat'l Conf. Artificial Intelligence Workshop Learning for Text Categorization*, vol. 752, 1998.
- [9] H. Schneiderman and T. Kanade, "A Statistical Approach to 3D Object Detection Applied to Faces and Cars," *Proc. IEEE CS Conf. Computer Vision and Pattern Recognition*, pp. 746-751, June 2000.
- [10] I. Kononenko, "A Semi-Naive Bayesian Classifier," *Proc. Sixth European Working Session on Learning*, pp. 206-219, Mar. 1991.
- [11] Y. Teh, M. Jordan, M. Beal, and D. Blei, "Hierarchical Dirichlet Processes," *J. Am. Statistical Assoc.*, vol. 101, no. 476, pp. 1566-1581, 2006.
- [12] D. Blei and M. Jordan, "Variational Inference for Dirichlet Process Mixtures," *J. Bayesian Analysis*, vol. 1, no. 1, pp. 121-144, 2005.
- [13] P. Muller and F. Quintana, "Nonparametric Bayesian Data Analysis," *Statistical Science*, vol. 19, no. 1, pp. 95-110, 2004.
- [14] F. Caron, M. Davy, A. Doucet, E. Duflos, and P. Vanheeghe, "Bayesian Inference for Linear Dynamic Models with Dirichlet Process Mixtures," *IEEE Trans. Signal Processing*, vol. 56, no. 1, pp. 71-84, Jan. 2008.
- [15] E. Jackson, M. Davy, A. Doucet, and W. Fitzgerald, "Bayesian Unsupervised Signal Classification by Dirichlet Process Mixtures of Gaussian Processes," *Proc. IEEE Int'l Conf. Acoustics, Speech, and Signal Processing*, Apr. 2007.
- [16] M. Davy, C. Doncarli, and J.Y. Tournet, "Classification of Chirp Signals Using Hierarchical Bayesian Learning and MCMC Methods," *IEEE Trans. Signal Processing*, vol. 50, no. 2, pp. 377-388, Feb. 2002.
- [17] C. Antoniak, "Mixtures of Dirichlet Processes with Applications to Bayesian Nonparametric Problems," *The Annals of Statistics*, vol. 2, pp. 1152-1174, 1974.
- [18] B. Shahbaba and R.M. Neal, "Nonlinear Models Using Dirichlet Process Mixtures," Technical Report 0707, Dept. of Statistics, Univ. of Toronto, 2007.
- [19] T. Ferguson, "A Bayesian Analysis of Some Nonparametric Problems," *The Annals of Statistics*, vol. 1, pp. 209-230, 1973.

- [20] J. Sethuraman, "A Constructive Definition of Dirichlet Priors," *Statistica Sinica*, vol. 4, pp. 639-650, 1994.
- [21] M. Escobar and M. West, "Bayesian Density Estimation and Inference Using Mixtures," *J. Am. Statistical Assoc.*, vol. 90, pp. 577-588, 1995.
- [22] D. Blackwell and J. MacQueen, "Ferguson Distributions via Polya Urn Schemes," *The Annals of Statistics*, vol. 1, no. 2, pp. 353-355, 1973.
- [23] S. Richardson and P.J. Green, "On Bayesian Analysis of Mixtures with an Unknown Number of Components," *J. Royal Statistical Soc. Series B*, vol. 59, no. 4, pp. 731-792, 1997.
- [24] J.-Y. Tournet, C. Mailhes, L. Amarouche, and N. Steunou, "Classification of Altimetric Signals Using Linear Discriminant Analysis," *Proc. IEEE Int'l Geoscience and Remote Sensing Symp.*, July 2008.
- [25] C.J. Oliver and S. Quegan, *Understanding Synthetic Aperture Radar Images*. Artech House, 1998.
- [26] J. Severini, C. Mailhes, P. Thibault, and J.-Y. Tournet, "Bayesian Estimation of Altimeter Echo Parameters," *Proc. IEEE Int'l Geoscience and Remote Sensing Symp.*, July 2008.
- [27] A. Gelman, J.B. Carlin, H.S. Stern, and D.B. Rubin, *Bayesian Data Analysis*, second ed., Chapman & Hall/CRC, 2004.
- [28] G. Brown, "The Average Impulse Response of a Rough Surface and Its Applications," *IEEE Trans. Antennas and Propagation*, vol. 25, no. 1, pp. 67-74, Jan. 1977.
- [29] S. Baker, O. Bombaci, C. Zeli, P. Venditti, O.-Z. Zanife, B. Soussi, J.-P. Dumont, J. Stum, M.P. Milagro-Perez, and J. Benveniste, "ENVISAT RA-2/MWR Product Handbook," [http://www.wcp.g.mssl.ucl.ac.uk/RA2\\_Handbook/concepts/ra2/ra2-mwr-PH-2.7.1.html#pgfid-868561](http://www.wcp.g.mssl.ucl.ac.uk/RA2_Handbook/concepts/ra2/ra2-mwr-PH-2.7.1.html#pgfid-868561), Mar. 2002.
- [30] J.P. Dumont, "ALT-RET-OCE-02—to Perform the Ocean-2 Retracking—Definition, Accuracy and Specification," technical report, Collecte Localisation Satellite (CLS), Toulouse, France, Oct. 2001.
- [31] C. Fraley and A.E. Raftery, "Bayesian Regularization for Normal Mixture Estimation and Model-Based Clustering," Technical Report 05/486, Dept. of Statistics, Univ. of Washington, Seattle, Aug. 2005.
- [32] T.W. Anderson, *An Introduction to Multivariate Statistical Analysis*, third ed. Wiley, 2003.
- [33] R. Herbrich, *Learning Kernel Classifiers*. MIT Press, 2002.
- [34] R. Neal, "Markov Chain Sampling Methods for Dirichlet Process Mixture Models," Technical Report 9815, Dept. of Statistics and Dept. of Computer Science, Univ. of Toronto, Ontario, Canada, 1998.



**Manuel Davy** received the engineer degree in electrical engineering from the Ecole Centrale de Nantes, France, in 1996, and the PhD degree from the University of Nantes in 2000. From 2000 to 2002, he was a research associate in the Signal Processing Group, University of Cambridge. In 2007, he created VEKIA in Lille, France, a company dedicated to software edition for forecasting solutions in retail. He is a member of the IEEE.



**Jean-Yves Tournet** received the ingénieur degree in electrical engineering from the Ecole Nationale Supérieure d'Electronique, d'Electrotechnique, d'Informatique et d'Hydraulique in Toulouse (ENSEEIH) in 1989 and the PhD degree from the National Polytechnic Institute, Toulouse in 1992. He is currently a professor at the University of Toulouse, France (ENSEEIH) and a member of the IRIT Laboratory (UMR 5505 of the CNRS). His research activities are centered around statistical signal processing, with a particular interest in classification and Markov Chain Monte Carlo methods. He was the program chair of the European Conference on Signal Processing (EUSIPCO), which was held in Toulouse in 2002. He was also member of the organizing committee for the international conference ICASSP '06 which was held in Toulouse in 2006. He has been a member of different technical committees. He is currently serving as an associate editor for the *IEEE Transactions on Signal Processing*. He is a senior member of the IEEE.