

# Classification bayésienne supervisée par processus de Dirichlet

Manuel DAVY<sup>1</sup>, Jean-Yves TOURNERET<sup>2</sup>

<sup>1</sup>CNRS/LAGIS et INRIA Futur équipe SequeL,  
40 Avenue Halley, 59650 Villeneuve d'Ascq, France

<sup>2</sup>ENSEEIH/Institut de Recherche en Informatique de Toulouse  
2, rue Charles Camichel, BP 7122 - F 31071 Toulouse Cedex 7, France  
Manuel.Davy@inria.fr, jean-yves.tourneret@enseeiht.fr

**Résumé** – La classification supervisée bayésienne de signaux, avec modèle génératif, comporte typiquement deux difficultés. La première est liée à la modélisation statistique des classes (c'est-à-dire à la définition de distributions *a priori* pour les paramètres). La seconde provient des difficultés liées au calcul bayésien. Dans cet article, nous proposons un *a priori* non paramétrique de type «mélange infini de distributions» connu sous le nom de *mélange de processus de Dirichlet*, ainsi qu'un algorithme de calcul bayésien par chaînes de Markov adapté. Ce couple modèle/algorithme est appliqué à la classification de données altimétriques radar.

**Abstract** – Generative bayesian supervised classification of signals is made difficult by two main issues. The first is related to the definition of a convenient prior distribution for the signal parameter. The second is in computing the quantities needed to perform the classification. In this paper, we propose a nonparametric prior known as the Dirichlet Process Mixture, as well as a suited Monte Carlo Markov Chain algorithm. These model and algorithm are applied to the classification of radar altimetry data.

## 1 Introduction

Supposons que l'on souhaite classer des signaux (données vectorielles) en  $K$  classes. Dans le contexte supervisé, on suppose en outre qu'un ensemble d'apprentissage est disponible, noté  $\mathbf{X}_j = \{\mathbf{x}_{1,j}, \dots, \mathbf{x}_{N_j,j}\}$  pour chaque classe  $j = 1, \dots, K$ . Pour chacune des classes, on suppose le modèle génératif des données suivant

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

où le paramètre vectoriel  $\theta_{i,j}$  appartient à l'espace  $\Theta_j$ ,  $i = 1, \dots, N_j$  et  $\epsilon_{i,j}$  est un bruit aléatoire, supposé additif ici à titre d'exemple. En termes statistiques,

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

où  $L_j$  est la vraisemblance spécifique à la classe numéro  $j$ , pour le paramètre  $\theta_{i,j}$ . Le but de l'apprentissage bayésien est d'apprendre la distribution de  $\theta_{i,j}$  à partir de l'ensemble d'apprentissage  $\mathbf{X}_j$ , pour chaque classe  $j = 1, \dots, K$ .

**Approche paramétrique hiérarchique.** L'approche paramétrique hiérarchique consiste à définir une distribution *a priori* sur l'espace du paramètre  $\Theta_j$ , notée  $g_j(\theta_j|\phi_j)$ , d'hyperparamètre  $\phi_j$ . La distribution *a posteriori* du paramètre est alors

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

où, par la règle de Bayes,  $p(\phi_j|\mathbf{X}_j) =$

$$\int \dots \int p(\theta_{1,j}, \dots, \theta_{N_j,j}, \phi_j|\mathbf{X}_j)d\theta_{1,j} \dots d\theta_{N_j,j} \quad (4)$$

où  $p(\theta_{1,j}, \dots, \theta_{N_j,j}, \phi_j|\mathbf{X}_j) =$

$$h_j(\phi_j) \prod_{i=1}^{N_j} [L_j(\mathbf{x}_{i,j}|\theta_{i,j})g_j(\theta_{i,j}|\phi_j)] \quad (5)$$

On voit que l'apprentissage de la distribution *a posteriori* du paramètre de chaque classe nécessite la définition d'une distribution *a priori* sur l'hyperparamètre, notée  $h_j(\phi_j)$ . Quand  $p(\theta_j|\mathbf{X}_j)$  est apprise, une nouvelle donnée  $\mathbf{x}$  est affectée à la classe de plus forte vraisemblance<sup>1</sup>  $p(\mathbf{x}|\mathbf{X}_j)$ ,  $j = 1, \dots, K$ , où

$$p(\mathbf{x}|\mathbf{X}_j) = \int L_j(\mathbf{x}|\theta_j)p(\theta_j|\mathbf{X}_j)d\theta_j \quad (6)$$

Cette approche peut être mise en œuvre à l'aide de l'algorithme MCMC présenté dans [4]. Un défaut de cette approche est de vouloir modéliser la distribution du paramètre par une distribution *a priori* paramétrique  $g_j(\theta_j|\phi_j)$ , dont le choix peut être extrêmement arbitraire dans les applications.

**Approche proposée.** Nous proposons une modélisation plus flexible, plus robuste et numériquement accessible, de la distribution du paramètre. Elle repose sur un modèle de mélange infini, dont la distribution de mélange est un processus de Dirichlet (DP).

## 2 Processus de Dirichlet et mélanges

Les DP [6, 1] peuvent être vus comme des distributions sur l'espace des distributions. Plus précisément, considérons une distribution de probabilité  $F_0$  définie sur un espace mesurable  $(\Theta, \Theta)$ . Une distribution  $F$  sur  $(\Theta, \Theta)$  est dite distribuée selon un DP (noté  $F \sim \mathcal{DP}(F; \alpha, F_0)$ ) si, pour  $A_1, \dots, A_m \in \Theta$  formant une partition de  $\Theta$ , on a

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

où  $\mathcal{D}$  est la distribution de Dirichlet habituelle. Les DP ont deux paramètres : la distribution de base  $F_0$  et le paramètre  $\alpha > 0$ .

<sup>1</sup>Nous supposons que toutes les classes ont la même probabilité *a priori*.

Une réalisation  $F$  d'un DP est presque sûrement discrète [1], et presque sûrement,

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

La représentation *stick breaking* fournit les poids  $\omega_k$  et les positions  $U_k$ , en répétant pour  $k = 1, 2, \dots$ , les étapes suivantes : 1) Echantillonner  $U_k \sim F_0(\cdot)$  et 2) Echantillonner  $\beta_k \sim \mathcal{B}(1, \alpha)$  puis calculer  $\omega_k = \beta_k \prod_{k'=1}^{k-1} (1 - \beta_{k'})$ , où  $\mathcal{B}$  est la distribution beta. On voit que  $F_0$  détermine les positions des composantes discrètes de  $F$ , tandis que  $\alpha$  règle la variance des poids  $\omega_k$ . Les DP sont particulièrement adaptés à l'échantillonnage conditionnel. En effet, considérons un ensemble de variables  $\{\psi_1, \dots, \psi_N\}$  dans  $\Psi$  i.i.d selon  $F(\cdot)$ , où  $F \sim \mathcal{DP}(F; \alpha, F_0)$ . Alors, pour  $i = 1, \dots, N$ ,

$$P(d\psi_i | \psi_{-i}, F_0, \alpha) = \frac{\alpha}{\alpha + N - 1} F_0(d\psi_i) + \frac{1}{\alpha + N - 1} \sum_{i'=1, \dots, N, i' \neq i} \delta_{\psi_{i'}}(d\psi_i) \quad (9)$$

où  $\psi_{-i} = \{\psi_{i'}\}_{i'=1, \dots, N, i' \neq i}$ . Cette formule est connue sous le nom d'*urne de Polya* [2].

**Mélanges par Processus de Dirichlet.** Une variable aléatoire  $\theta$  est distribuée selon un mélange par DP (DPM) quand elle est générée ainsi :

1. Echantillonner  $F \sim \mathcal{DP}(F; \alpha, F_0)$
2. Echantillonner  $\psi \sim F(\cdot)$
3. Echantillonner  $\theta \sim f(\theta | \psi)$

La répétition des étapes 2) et 3) génère une famille de variables aléatoires distribuées selon le DPM de densités  $f(\cdot | \cdot)$ . Chaque itération de l'étape 2) génère un paramètre  $\psi_i$  qui détermine le *cluster* auquel  $\theta_i$  est affecté. Considérant l'urne de Polya, nous remarquons que plusieurs  $\psi_i$  partagent la même position dans l'espace des paramètres. Dans la représentation *stick-breaking*, nous voyons que les positions sont  $U_k$  avec la probabilité  $\omega_k$ . Comme pour les mélanges finis, il est possible d'introduire des variables latentes  $z_i$  pour chaque  $\theta_i$  ( $i = 1, \dots, N$ ). Cela consiste à remplacer l'étape 2) ci-dessus par

- 2.1) Echantillonner  $z \sim P(z|F)$  où  $P(z = k|F) = \omega_k$ . Le poids  $\omega_k$  vient de la représentation *stick-breaking* de  $F$ .
- 2.2) Affecter  $\psi \leftarrow U_z$ , où  $U_z$  vient de la représentation *stick-breaking* de  $F$ .

Enfin, en notant le DPM par  $G(d\theta)$  on a :

$$G(d\theta) = \int_{\Theta} F(d\theta | \psi) dF(\psi) \text{ with } F \sim \mathcal{DP}(F; F_0, \alpha) \quad (10)$$

où  $F(\cdot | \psi)$  est la distribution de densité  $f(\cdot | \psi)$  sous la mesure dominante.

### 3 Modéliser la distribution du paramètre par un DPM

Dans cette section, nous remplaçons la distribution *a priori* du paramètre  $\int_{\Phi} g(\theta | \phi) h(\phi) d\phi$ , introduite à la section 1 par le DPM  $G(d\theta)$ . Par soucis de clarté des notations, l'indice de

classe est omis. Des calculs similaires sont effectués indépendamment dans chaque classe. Nous considérons le modèle hiérarchique, pour chaque classe

$$F \sim \mathcal{DP}(F; F_0, \alpha) \quad (11)$$

$$\psi_i \sim F \quad (12)$$

$$\theta_i \sim f(\theta_i | \psi_i) \quad (13)$$

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

Il est possible d'écrire la distribution *a posteriori* du paramètre (en introduisant les variables latentes  $\mathbf{z} = [z_1, \dots, z_n]$  et les positions des *clusters*  $\mathbf{U} = \{U_1, U_2, \dots\}$ ),

$$\begin{aligned} p(\theta | \mathbf{X}) &= \int \sum_{\mathbf{z}} p(\theta, \mathbf{z}, \mathbf{U} | \mathbf{X}) d\mathbf{U} \\ &= \int \sum_{\mathbf{z}, \mathbf{z}} f(\theta | U_z) P(\mathbf{z} | \mathbf{z}) P(d\mathbf{U}, \mathbf{z} | \mathbf{X}) \end{aligned} \quad (15)$$

où  $P(\mathbf{z} | \mathbf{z})$  vérifie  $P(z = k | \mathbf{z}) = \frac{1}{N} \sum_{i=1}^N \delta_{k, z_i}$ . La distribution *a posteriori* des paramètres suit

$$\begin{aligned} P(d\mathbf{U}, \mathbf{z} | \mathbf{X}) &= \int_{\Theta^N} P(d\theta, d\mathbf{U}, \mathbf{z} | \mathbf{X}) \\ &\propto \int_{\Theta^N} \mathcal{DP}(F; F_0, \alpha) \prod_{i=1}^N \mathbf{L}(\mathbf{x}_i | \theta_i) f(\theta_i | U_{z_i}) \end{aligned} \quad (16)$$

L'inférence est effectuée en appliquant l'algorithme ci-dessous.

---

#### Algorithme 1: Apprentissage bayésien par MCMC

---

% Step 0 : Initialisation

– Pour  $i = 1, \dots, N$ , échantillonner  $\tilde{\psi}^{(0)} \sim p(\psi | \tilde{\psi}_1^{(0)}, \dots, \tilde{\psi}_{i-1}^{(0)})$  par l'urne de Polya et en déduire  $\tilde{\mathbf{z}}^{(0)}$  et  $\tilde{\mathbf{U}}^{(0)}$

– Pour  $i = 1, \dots, N$ , échantillonner  $\tilde{\theta}_i^{(0)} \sim f(\theta | \tilde{\psi}_i^{(0)})$

% Step 1 : Itérations Pour  $l = 1, \dots, L$ , faire

1.1- Échantillonner  $\tilde{\mathbf{z}}^{(l)} \sim P(\mathbf{z} | \tilde{\mathbf{U}}^{(l-1)}, \tilde{\theta}^{(l-1)})$  directement (dans le cas où  $F_0$  et  $f(\cdot | \cdot)$  sont conjuguée) ou par un pas de Métropolis-Hastings (MH)

1.2- Échantillonner  $\tilde{\mathbf{U}}^{(l)} \sim p(\mathbf{U} | \tilde{\mathbf{z}}^{(l)}, \tilde{\theta}^{(l-1)})$ , directement ou par un pas de MH

1.3- Échantillonner  $\tilde{\theta}^{(l)} \sim p(\theta | \tilde{\mathbf{z}}^{(l)}, \tilde{\mathbf{U}}^{(l)})$  : pour  $i = 1, \dots, N$ , échantillonner par un pas de MH  $\tilde{\theta}_i^{(l)} \sim p(\theta_i | \mathbf{x}_i, \tilde{U}_{\tilde{z}_i^{(l)}}^{(l)}) \propto$

$$\mathbf{L}(\mathbf{x}_i | \theta_i) f(\theta_i | \tilde{U}_{\tilde{z}_i^{(l)}}^{(l)})$$

1.4- Échantillonner  $\tilde{\xi}^{(l)} \sim P(\xi | \tilde{\mathbf{z}}^{(l)})$  tel que  $P(\xi = k | \tilde{\mathbf{z}}^{(l)}) = \frac{1}{N} \sum_{i=1}^N \delta_{k, \tilde{z}_i^{(l)}}$

1.5- Échantillonner  $\tilde{\theta}^{(l)} \sim f(\theta | \tilde{\xi}^{(l)})$

1.6- (Optionel) Echantillonner les hyperparamètres de  $F_0$  et  $\alpha$

---

Les étapes 1.1 et 1.2 de l'algorithme ci-dessus sont effectuées selon l'une des procédures proposées dans [8], et l'hyperparamètre  $\alpha$  peut être mis à jour selon la technique de [9]. Une remarque d'une grande utilité pratique est que le processus de Dirichlet est représenté sous la forme de centroïdes de clusters  $\mathbf{U}$  et de variables latentes  $\mathbf{z}$ . Cette paramétrisation fait qu'aucune troncature de la représentation *stick breaking* n'est réalisée explicitement : seuls les *clusters* actifs sont mis à jour dans l'algorithme (le *clusters*  $j$  est actif si au moins un paramètre  $\theta_i$  y est associé par la variable latente  $z_i = j$ ).

## 4 Application à la classification de données altimétriques radar

Le satellite Poseïdon est équipé d'un radar altimétrique fournissant des données sous forme de séries temporelles de longueur utile  $T = 104$  échantillons, dont 10 sont représentées à la figure 1. Elles peuvent être décrites par le modèle de Hayne [3, 7, 5], valable pour les zones d'océans, avec

$$\mathbf{x}_t = F(t, \theta) \epsilon_t, \quad t = 1, \dots, T \quad (17)$$

où le bruit  $\epsilon_t$  est blanc de distribution gamma  $\mathcal{G}a(100, 100)$  et

$$F(t, \theta) = P_n + \frac{a_c \sigma_0}{2} \left[ 1 + \operatorname{erf} \left( \frac{t - \tau - c_c \sigma_c^2}{\sqrt{2} \sigma_c} \right) \right] \exp \left[ -c_c \left( t - \tau - \frac{c_c \sigma_c^2}{2} \right) \right] \quad (18)$$

avec

$$a_c = \exp \left( \frac{-4 \sin^2 \zeta}{\xi_2} \right) \quad (19)$$

$$c_c = \xi_1 \left[ \cos(2\zeta) - \frac{\sin^2(2\zeta)}{\xi_2} \right] \quad (20)$$

Le vecteur paramètre inconnu est

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

avec  $\tau > 0$  le temps de propagation de l'onde radar,  $\sigma_0 > 0$  un coefficient de rétrodiffusion spécifique à la zone d'océan observée,  $\text{SWH} = 2c_{\text{light}} \sigma_s$  où  $\sigma_s^2 = \sigma_c^2 - \sigma_p^2$  ( $c_{\text{light}}$  est la vitesse de la lumière et  $\sigma_p^2$  est connu) est la hauteur des vagues et  $\zeta$  est l'angle de dépointage de l'antenne. Toutes les autres grandeurs sont connues. La vraisemblance est donc

$$L(\mathbf{x}|\theta) \propto \exp \left( -L \sum_{t=1}^T \left[ \log F(t, \theta) + \frac{\mathbf{x}_t}{F(t, \theta)} \right] \right) \quad (21)$$

La distribution de mélange est supposée gaussienne  $f(\theta|\psi) = \mathcal{N}(\mu, \Sigma)$  où  $\psi = \{\mu, \Sigma\}$  contient le vecteur moyenne  $\mu$  et la matrice de covariance  $\Sigma$ . Par simplicité, nous supposons que  $F_0$  est une normale inverse Wishart

$$F_0(\psi) = \mathcal{N}(\mu; \mu_0, \Sigma/\kappa_0) \mathcal{IW}(\Sigma; \nu_0, \Lambda_0), \quad (22)$$

où  $\mu_0, \kappa_0, \nu_0, \Lambda_0$  sont des hyperparamètres.

L'algorithme 1 est mis en œuvre ici, dans une version où les hyperparamètres  $\kappa_0, \nu_0$  et  $\Lambda_0$  sont pourvus de distributions *a priori* et échantillonnés dans l'étape 1.6 de l'algorithme 1.

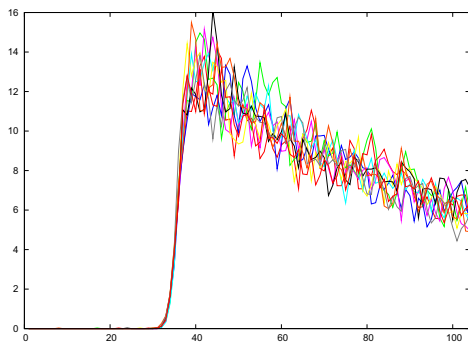


FIG. 1 – Dix signaux radar altimétriques obtenus au dessus d'océans.

Afin de vérifier la capacité de l'algorithme et du modèle à bien apprendre la distribution *a posterior* du paramètre, nous avons effectué 100 exécutions de l'algorithme 1 pour  $N = 100$  signaux d'apprentissage sur 5000 itérations, dont 3000 itérations de chauffage. Pour chaque exécution, les estimées  $\hat{\theta}_i^{\text{MMSE}}$  ont été comparées aux valeurs  $\theta_i^{\text{gt}}$  utilisées pour synthétiser les données (pour  $i = 1, \dots, 100$ ) via l'indice d'erreur quadratique

$$\text{Err}_{\theta}(N) = \frac{1}{N} \sum_{i=1}^N (\hat{\theta}_i^{\text{MMSE}} - \theta_i^{\text{gt}})^{\top} \mathbf{S}^{-1} (\hat{\theta}_i^{\text{MMSE}} - \theta_i^{\text{gt}}) \quad (23)$$

où  $\mathbf{S}^{-1} = \text{diag}(\tau^{\text{ref}}, \sigma_0^{\text{ref}}, \text{SWH}^{\text{ref}}, \zeta^{\text{ref}})$  et

$$\hat{\theta}_i^{\text{MMSE}} = (L - L_{\text{burn in}})^{-1} \sum_{l=L_{\text{burn in}}+1}^L \tilde{\theta}_i^{(l)} \quad (24)$$

La moyenne sur 100 exécutions de  $\text{Err}_{\theta}(N)$  est de 0.109 avec un écart type de 0.063. La figure 2 montre l'évolution moyenne de l'estimation au cours des itérations.

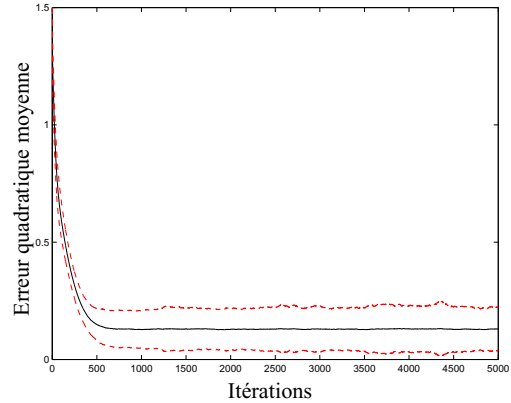


FIG. 2 – Evolution de l'erreur quadratique moyenne du paramètre  $\theta_i, i = 1, \dots, 100$  au cours des 5000 itérations de l'algorithme MCMC, moyenné sur 100 exécutions. Les courbes pointillées représentent l'intervalle de 95% de confiance ( $\pm 2$  écart-types).

La figure 3 représente les histogrammes des quatre composantes des échantillons  $\tilde{\theta}^{(l)}, l \in \{3001, \dots, 5000\}$  pour une exécution typique de l'algorithme 1, ainsi que les densités de probabilité utilisées pour générer les valeurs des paramètres utilisés en simulation. Dans le cas du troisième paramètre, le mélange de processus de Dirichlet a permis de bien apprendre les deux modes de la distribution. La figure 4 représente le même type de résultats, dans le cas où seuls  $N = 10$  signaux d'apprentissage sont utilisés. Là encore, on voit que le modèle de DPM apprend bien la distribution du paramètre. L'histogramme de la composante 3 n'est pas bimodal, ce qui montre que le modèle régularise bien lorsque le nombre de données d'apprentissage est faible. Mettre deux modes serait revenu à faire une hypothèse forte compte tenu de la faible information disponible. Enfin, la figure 5 représente ces histogrammes dans le cas où  $N = 100$  données d'apprentissage réelles sont utilisées.

## 5 Discussion, Conclusion et Perspectives

La modélisation proposée permet une grande flexibilité tout en contrôlant la complexité du modèle. Les simulations ont

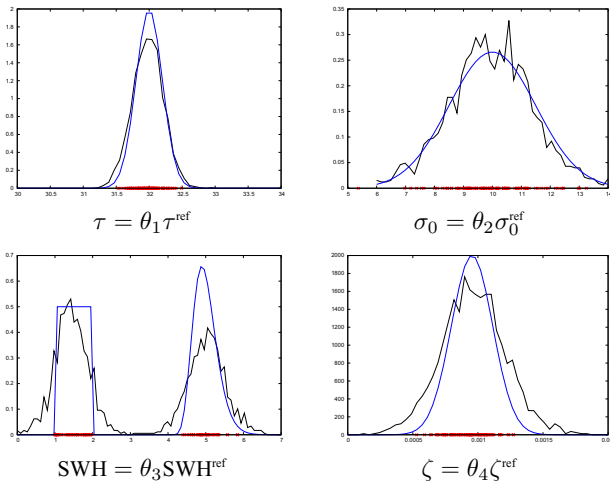


FIG. 3 – En noir, histogrammes des quatre composantes du vecteur paramètre calculés à partir des échantillons  $\hat{\theta}^{(l)}$  pour  $l \in \{3001, \dots, 5000\}$ . En bleu, densités de probabilité utilisées pour tirer les quatre paramètres (représentés par les croix rouges) lors de la synthèse des  $N = 100$  signaux.

montré qu'il était possible d'apprendre des distributions multimodales ou plus généralement, non gaussiennes. Le choix de la distribution de base  $F_0$  est peu contraignant en pratique, en particulier du fait que ses paramètres sont estimés aussi. Les résultats de simulation présentés montrent la capacité d'apprentissage du modèle, dont découle directement l'efficacité de classification. L'algorithme MCMC proposé est assez simple à mettre en œuvre, et il peut s'écrire sous forme générique pour la plupart des étapes.

La suite de ce travail consiste à développer un algorithme de type Monte Carlo Séquentiel pour accélérer la convergence de l'algorithme MCMC présenté ici, et de traiter d'autres applications de classification supervisée où un modèle génératif est disponible.

## 6 Remerciements

Les auteurs remercient vivement la société CLS et plus particulièrement Laiba Amarouche, Pierre Thibault et Ouan-Zan Zanife pour leurs nombreuses discussions concernant l'altimétrie radar.

## Références

- [1] C.E. Antoniak. Mixtures of Dirichlet processes with applications to Bayesian nonparametric problems. *The annals of statistics*, 2 :1152–1174, 1974.
- [2] D. Blackwell and J.B. MacQueen. Ferguson distributions via Polya urn schemes. *The Annals of Statistics*, 1(2) :353–355, 1973.
- [3] G. Brown. The average impulse response of a rough surface and its applications. *IEEE Transactions on Antennas and Propagation*, 25(1) :67–74, 1977.
- [4] M. Davy, C. Doncarli, and J. Y. Tourneret. Classification of chirp signals using hierarchical Bayesian learning and MCMC methods. *IEEE transactions on Signal Processing*, 50(2) :377–388, February 2002.

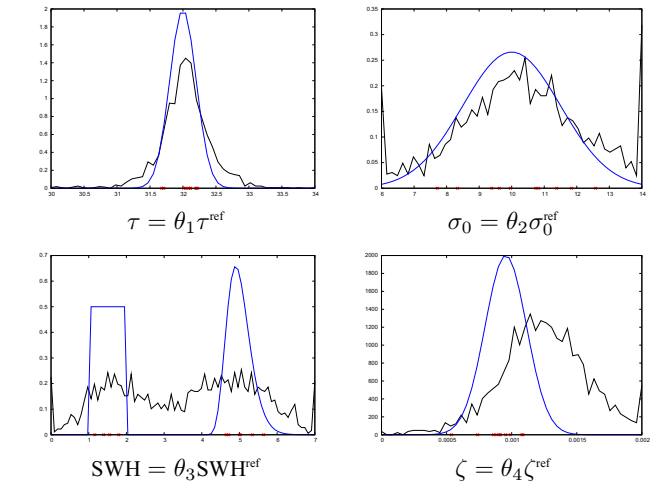


FIG. 4 – Mêmes histogrammes qu'à la figure 3 dans le cas où le nombre de données d'apprentissage est  $N = 10$

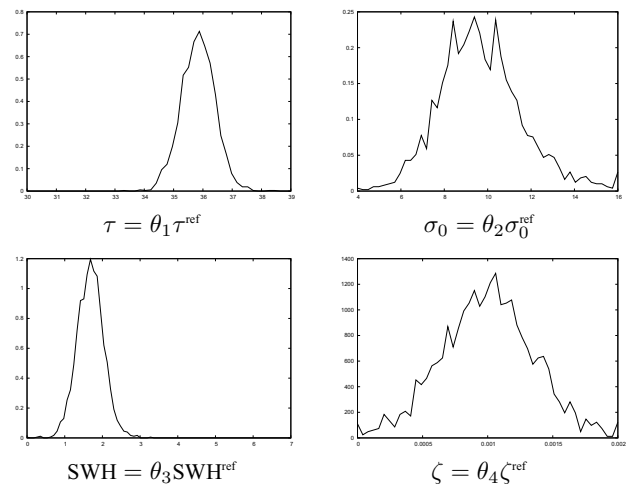


FIG. 5 – Mêmes histogrammes qu'à la figure 3 dans le cas où l'ensemble d'apprentissage est composé de  $N = 100$  données réelles.

- [5] 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, October 2001.
- [6] T.S. Ferguson. A Bayesian analysis of some nonparametric problems. *The annals of statistics*, 1 :209–230, 1973.
- [7] G. Hayne. Radar altimeter mean return waveforms from near-normal-incidence ocean surface scattering. *IEEE Transactions on Antennas and Propagation*, 28(5) :687–692, 1980.
- [8] R.M. Neal. Markov chain sampling methods for Dirichlet process mixture models. Technical Report 9815, Dpt of Statistics and department of computer science, University of Toronto, Ontario, Canada, 1998.
- [9] M. West. Hyperparameter estimation in Dirichlet process mixture models. Technical report, Institute of Statistics and Decision Sciences, Duke university, 1992.