

Apprentissage statistique à l'ère du *Big Data* :  
stabilité et efficacité d'approximations de certains  
algorithmes

Florian Maire

Lundi 8 janvier 2018

# Programme de cette présentation

Quelques informations sur mon parcours

Statistique computationnelle : modélisation et apprentissage

Apprentissage statistique en présence d'un grand nombre de données : stabilité et efficacité d'approximations de certains algorithmes.

Quelques Illustrations

# Plan

Quelques informations sur mon parcours

Statistique computationnelle : modélisation et apprentissage

Apprentissage statistique en présence d'un grand nombre de données : stabilité et efficacité d'approximations de certains algorithmes.

Quelques Illustrations

## Quelques dates

- ▶ 2005–2007, Classes prépa, Lycée Lavoisier, Paris
- ▶ 2007–2010, École d'ingénieur, Telecom Sudparis, Evry



- ▶ 2011–2014, Doctorat, Université Paris VI, Telecom Sudparis et ONERA



- ▶ 2014–2017, Post-doctorat, University College Dublin, Irlande



**THESE DE DOCTORAT TELECOM SUDPARIS et  
UNIVERSITE PIERRE ET MARIE CURIE**

**Spécialité : Probabilités et Statistiques**

**Ecole doctorale : Informatique, Télécommunications et  
Electronique de Paris**

**Présentée par**

**Florian MAIRE**

**Pour obtenir le grade de  
DOCTEUR DE TELECOM SUDPARIS**

**DETECTION ET CLASSIFICATION DE CIBLES  
MULTISPECTRALES DANS L'INFRAROUGE**

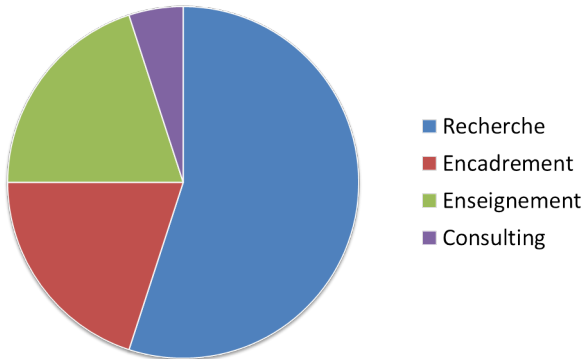
dirigée par Randal DOUC, Professeur à Telecom SudParis,  
co-dirigée par Eric MOULINES, Professeur à Telecom ParisTech,  
encadrée par Sidonie LEFEBVRE, Ingénieur de Recherche à l'ONERA.

**soutenue le : 14 Février 2014**

**devant le jury composé de :**

M <sup>me</sup> Stéphanie ALLASSONNIÈRE	Ecole Polytechnique	examinateur
M. Jacques BLANC-TALON	DGA	examinateur
M. Randal DOUC	Telecom SudParis	directeur
M <sup>me</sup> Sidonie LEFEBVRE	ONERA	encadrant
M. Jean-Michel MARIN	Université de Montpellier 2	rapporteur
M. Eric MOULINES	Telecom ParisTech	co-directeur
M. Christian ROBERT	Université Paris-Dauphine	examinateur
M. Olivier WINTENBERGER	Université Paris 6	examinateur

## Post-doctorat



- ▶ Recherche : Apprentissage en statistique Bayésienne (groupe de Prof. Nial Friel)
- ▶ Encadrement : PhD de Lampros Bourannis (avec N. Friel), 2 thèses de Masters et 3 étudiants de Masters
- ▶ Enseignement : Probability Theory, Time Series, Monte Carlo methods

# Plan

Quelques informations sur mon parcours

Statistique computationnelle : modélisation et apprentissage

Apprentissage statistique en présence d'un grand nombre de données : stabilité et efficacité d'approximations de certains algorithmes.

Quelques Illustrations

# Notations

Les **données**  $\{Y_n, n \in \mathbb{N}\}$  sont des observations d'un phénomène aléatoires  $\mathcal{Y}$  d'intérêt

$$\mathcal{Y} \rightsquigarrow Y_1, Y_2, \dots$$

**Exemples** : signal, image, graphe

- ▶ les données sont généralement dans un espace euclidien de dimension  $d$

$$Y_1 \in \mathcal{Y} \subseteq \mathbb{R}^d.$$

- ▶ Un **modèle** est une distribution de probabilité  $f$  sur l'espace mesurable  $(\mathcal{Y}, \mathcal{Y})$
- ▶  $f \in \mathcal{F}$  où  $\mathcal{F}$  est une famille de distributions de proba. sur  $(\mathcal{Y}, \mathcal{Y})$

**Apprentissage** : trouver un "bon" modèle  $f^*$  appartenant à la famille  $\mathcal{F}$

ayant observé  $Y = Y_1, Y_2, \dots, Y_N$

$$\Pr(Y \in A) \approx f^*(A) \quad \forall A \in \mathcal{Y}$$



# Statistique paramétrique

**Hypothèse** :  $\mathcal{F}$  est une famille de distributions paramétrisée par un vecteur  $\theta$

$$\theta \in \Theta \subseteq \mathbb{R}^p \quad (p < \infty)$$

$$\Rightarrow \forall f \in \mathcal{F}, \exists \theta \in \Theta, f \equiv f(\theta)$$

$\Rightarrow$  Passage à un problème en dimension fini

**Apprentissage** : trouver un "bon" vecteur  $\theta^* \in \Theta$  de sorte à ce qu'ayant observé

$$Y = Y_1, Y_2, \dots, Y_N$$

$$\Rightarrow \Pr(Y \in A) \approx f(A | \theta^*), \quad A \in \mathcal{Y}$$

## Exemple 1 : Regression logistique

**Application** : Détecter la présence d'une maladie chez un patient

- ▶ Les données sont des réponses booléennes  $Y_1 \in \{0, 1\}$   
*patient 1 est malade ou sain*
- ▶ Variables explicatives (covariates)  $X_1 = (X_{1,1}, \dots, X_{1,r})$   
*patient 1 a des données auxiliaires : age, sexe, indice de santé,...*

Le modèle s'écrit

$$f(Y_1 | X_1, \theta) = \left\{ \frac{1}{1 + \exp(-\theta^T X_1)} \right\}^{Y_1} \left\{ 1 - \frac{1}{1 + \exp(-\theta^T X_1)} \right\}^{1 - Y_1} .$$

**Apprentissage** : consiste à estimer le paramètre  $\theta \in \Theta$   
sachant  $\{(Y_1, X_1), \dots, (Y_n, X_n)\}$

## Exemple 2 : Modèle à templates déformables

### Application : Classifier un chiffre manuscrit

- ▶ l'image d'un chiffre  $\mathcal{Y}_n$  est modélisée par une collection de fonctions **déterministes** (les **templates**)  $\mathcal{T}^{(0)}, \dots, \mathcal{T}^{(9)}$ :

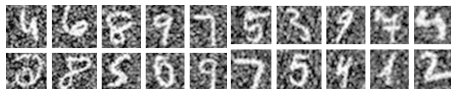
$$\mathcal{T}^{(i)} : \mathbb{R}^2 \rightarrow \mathbb{R}$$

- ▶ ces templates sont observés à travers un processus aléatoire
  - (i) de déformation du plan  $D_n : \mathbb{R}^2 \rightarrow \mathbb{R}^2$ ,
  - (ii) et un sous bruit additif  $\mathcal{W}_n : \mathbb{R}^2 \rightarrow \mathbb{R}$
- ▶ Modèle:

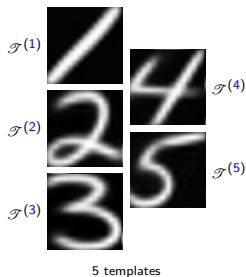
$$\forall u \in \mathbb{R}^2, \quad \mathcal{Y}_n(u) = \mathcal{T}^{(i)} \circ D_n(u) + \mathcal{W}_n(u).$$

- ▶ Les données sont des fonctions  $\{\mathcal{Y}_n : \mathbb{R}^2 \rightarrow \mathbb{R}, n \leq N\}$  discrétisées :

$$\{\mathcal{Y}_n, n \leq N\} \xrightarrow[\text{sur une grille de pixels}]{\text{discretisation}} \{Y_n, n \leq N\}$$



## Exemple 2 : Modèle à templates déformables



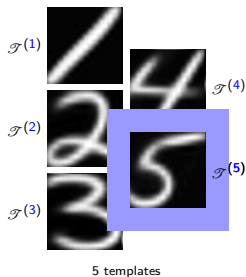
En paramétrisant

- ▶ la famille de templates  $\mathcal{T}$  par  $\theta \in \Theta$
- ▶ le champ de déformation  $D_n$  par  $\beta_n \in \mathbb{B}$

et en discrétisant  $\mathcal{Y}_n$  sur une grille, le modèle peut se réécrire comme :

$$Y_n = \Phi_{\beta_n} \theta_{i_n} + \zeta_n, \quad \zeta_n \sim \mathcal{N}(0, \Sigma).$$

## Exemple 2 : Modèle à templates déformables



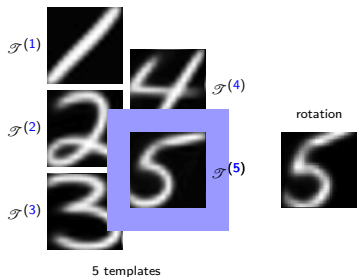
En paramétrisant

- ▶ la famille de templates  $\mathcal{T}$  par  $\theta \in \Theta$
- ▶ le champ de déformation  $D_n$  par  $\beta_n \in \mathbb{B}$

et en discrétisant  $\mathcal{Y}_n$  sur une grille, le modèle peut se réécrire comme :

$$Y_n = \Phi_{\beta_n} \theta_{i_n} + \zeta_n, \quad \zeta_n \sim \mathcal{N}(0, \Sigma).$$

## Exemple 2 : Modèle à templates déformables



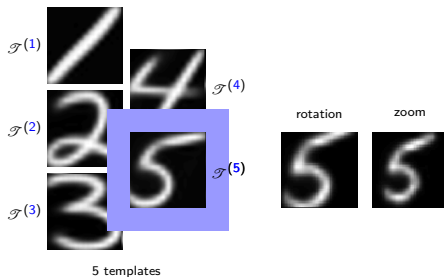
En paramétrisant

- ▶ la famille de templates  $\mathcal{T}$  par  $\theta \in \Theta$
- ▶ le champ de déformation  $D_n$  par  $\beta_n \in \mathbb{B}$

et en discrétisant  $\mathcal{Y}_n$  sur une grille, le modèle peut se réécrire comme :

$$Y_n = \Phi_{\beta_n} \theta_{i_n} + \zeta_n, \quad \zeta_n \sim \mathcal{N}(0, \Sigma).$$

## Exemple 2 : Modèle à templates déformables



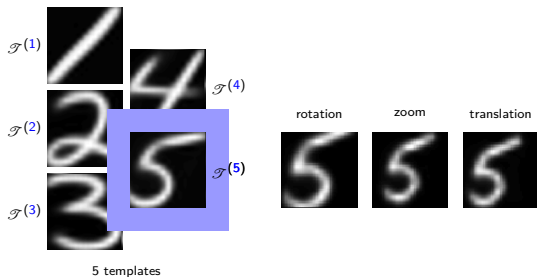
En paramétrisant

- ▶ la famille de templates  $\mathcal{T}$  par  $\theta \in \Theta$
- ▶ le champ de déformation  $D_n$  par  $\beta_n \in \mathbb{B}$

et en discrétisant  $\mathcal{Y}_n$  sur une grille, le modèle peut se réécrire comme :

$$Y_n = \Phi_{\beta_n} \theta_{i_n} + \zeta_n, \quad \zeta_n \sim \mathcal{N}(0, \Sigma).$$

## Exemple 2 : Modèle à templates déformables



En paramétrisant

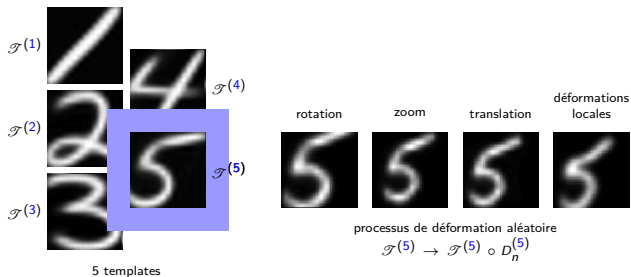
- ▶ la famille de templates  $\mathcal{T}$  par  $\theta \in \Theta$
- ▶ le champ de déformation  $D_n$  par  $\beta_n \in \mathbb{B}$

et en discrétisant  $\mathcal{Y}_n$  sur une grille, le modèle peut se réécrire comme :

$$Y_n = \Phi_{\beta_n} \theta_{i_n} + \zeta_n, \quad \zeta_n \sim \mathcal{N}(0, \Sigma).$$



## Exemple 2 : Modèle à templates déformables



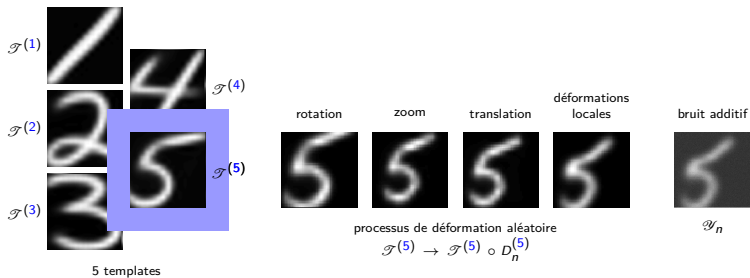
En paramétrisant

- ▶ la famille de templates  $\mathcal{T}$  par  $\theta \in \Theta$
- ▶ le champ de déformation  $D_n$  par  $\beta_n \in \mathbb{B}$

et en discrétisant  $\mathcal{Y}_n$  sur une grille, le modèle peut se réécrire comme :

$$Y_n = \Phi_{\beta_n} \theta_{i_n} + \zeta_n, \quad \zeta_n \sim \mathcal{N}(0, \Sigma).$$

## Exemple 2 : Modèle à templates déformables



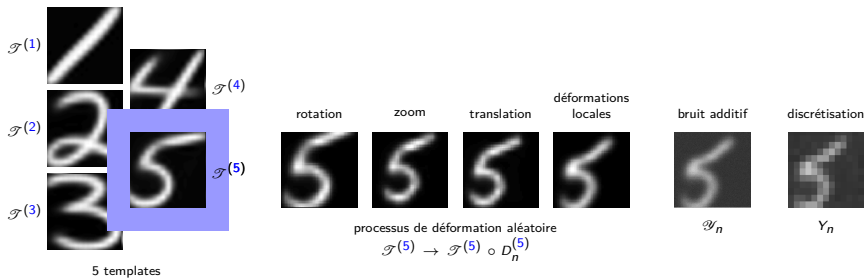
En paramétrisant

- ▶ la famille de templates  $\mathcal{T}$  par  $\theta \in \Theta$
- ▶ le champ de déformation  $D_n$  par  $\beta_n \in \mathbb{B}$

et en discrétisant  $\mathcal{Y}_n$  sur une grille, le modèle peut se réécrire comme :

$$Y_n = \Phi_{\beta_n} \theta_{i_n} + \zeta_n, \quad \zeta_n \sim \mathcal{N}(0, \Sigma).$$

## Exemple 2 : Modèle à templates déformables



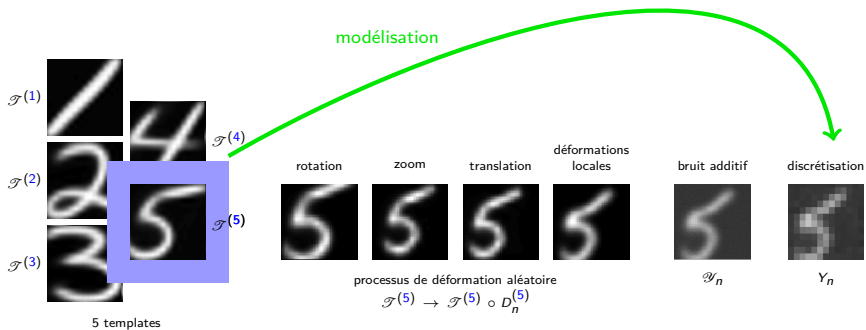
En paramétrisant

- ▶ la famille de templates  $\mathcal{T}$  par  $\theta \in \Theta$
- ▶ le champ de déformation  $D_n$  par  $\beta_n \in \mathbb{B}$

et en discrétisant  $\mathcal{Y}_n$  sur une grille, le modèle peut se réécrire comme :

$$Y_n = \Phi_{\beta_n} \theta_{i_n} + \zeta_n, \quad \zeta_n \sim \mathcal{N}(0, \Sigma).$$

## Exemple 2 : Modèle à templates déformables



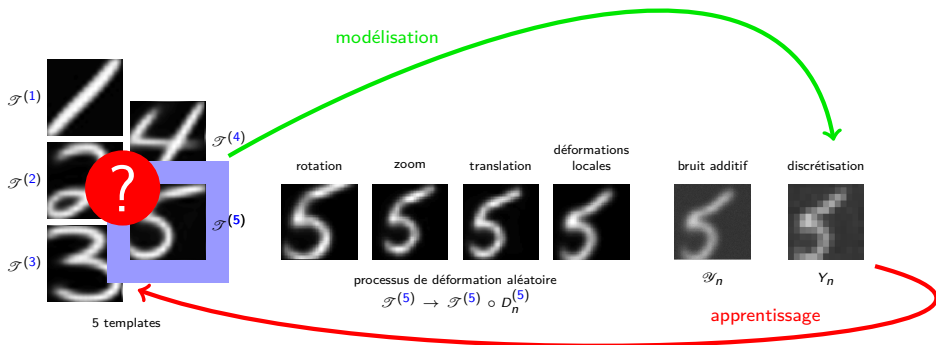
En paramétrisant

- ▶ la famille de templates  $\mathcal{T}$  par  $\theta \in \Theta$
- ▶ le champ de déformation  $D_n$  par  $\beta_n \in \mathbb{B}$

et en discrétisant  $\mathcal{Y}_n$  sur une grille, le modèle peut se réécrire comme :

$$Y_n = \Phi_{\beta_n} \theta_{i_n} + \zeta_n, \quad \zeta_n \sim \mathcal{N}(0, \Sigma).$$

## Exemple 2 : Modèle à templates déformables



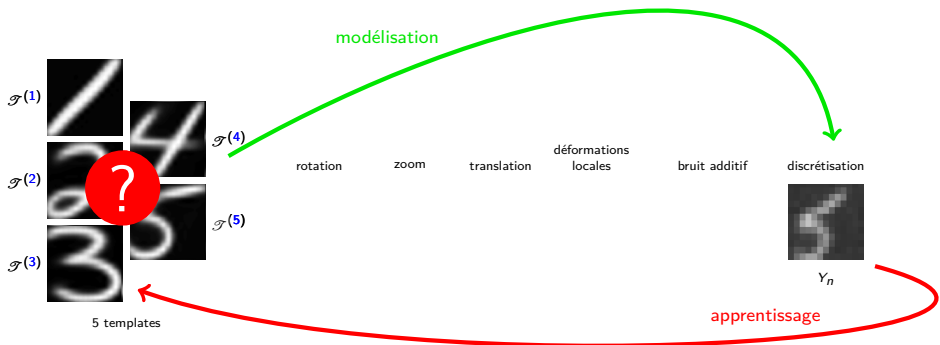
En paramétrisant

- ▶ la famille de templates  $\mathcal{T}$  par  $\theta \in \Theta$
- ▶ le champ de déformation  $D_n$  par  $\beta_n \in \mathbb{B}$

et en discrétisant  $\mathcal{Y}_n$  sur une grille, le modèle peut se réécrire comme :

$$Y_n = \Phi_{\beta_n} \theta_{i_n} + \zeta_n, \quad \zeta_n \sim \mathcal{N}(0, \Sigma).$$

## Exemple 2 : Modèle à templates déformables



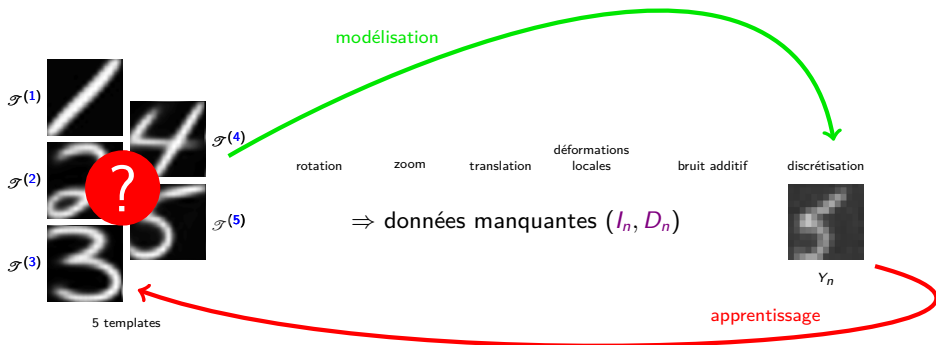
En paramétrisant

- ▶ la famille de templates  $\mathcal{T}$  par  $\theta \in \Theta$
- ▶ le champ de déformation  $D_n$  par  $\beta_n \in \mathbb{B}$

et en discrétisant  $\mathcal{Y}_n$  sur une grille, le modèle peut se réécrire comme :

$$Y_n = \Phi_{\beta_n} \theta_{i_n} + \zeta_n, \quad \zeta_n \sim \mathcal{N}(0, \Sigma).$$

## Exemple 2 : Modèle à templates déformables



En paramétrisant

- ▶ la famille de templates  $\mathcal{T}$  par  $\theta \in \Theta$
- ▶ le champ de déformation  $D_n$  par  $\beta_n \in \mathbb{B}$

et en discrétisant  $\mathcal{Y}_n$  sur une grille, le modèle peut se réécrire comme :

$$Y_n = \Phi_{\beta_n} \theta_{i_n} + \zeta_n, \quad \zeta_n \sim \mathcal{N}(0, \Sigma).$$

# Apprentissage paramétrique

## Deux écoles

	fréquentiste	Bayésienne
modèle	$f(\cdot   \theta), \theta \in \Theta$	$f(\cdot   \theta), \theta \in \Theta$
source d'information	données $Y = Y_1, Y_2, \dots$	données $Y = Y_1, Y_2, \dots$ + une distribution <i>a priori</i> $\theta \sim p$
apprentissage	trouver l'ensemble $\Theta^*$ $\arg \min_{\theta \in \Theta} d(Y, f(\cdot   \theta))$ pour une distance $d$	trouver la distribution $\pi$ $\pi(\theta   Y) \propto f(Y   \theta)p(\theta)$
nature du problème	optimisation	échantillonnage
objectif	$\Pr(\theta \in A   Y)$ intervalle de confiance	$\Pr(\theta \in A) = \pi(A   Y)$ intervalle de crédibilité
exemple	Maximum de vraisemblance	méthodes de Monte Carlo



## Exemple 1 (régression logistique) : apprentissage Bayésien

- ▶  $N = 10^7$  données  $Y = (Y_1, \dots, Y_N)$
- ▶ covariates  $X = (X_1, \dots, X_N)$
- ▶ loi a priori  $p(\theta) = \text{Laplace}(\theta; 1) \propto \exp\{-\|\theta\|_1\}$
- ▶ loi a posteriori

$$\pi(\theta | Y, X) \propto \prod_{n=1}^N \left\{ \frac{1}{1 + \exp(-\theta^T X_n)} \right\}^{Y_n} \left\{ 1 - \frac{1}{1 + \exp(-\theta^T X_n)} \right\}^{1 - Y_n} e^{-\|\theta\|_1}$$

- ▶  $\pi$  n'est pas une distribution standard  $\Rightarrow$  l'échantillonnage de  $\pi$  n'est pas simple

**Apprentissage Bayésien** : les méthodes de Monte Carlo par chaîne de Markov (MCMC) offrent une possibilité simple et universelle d'échantillonner  $\pi$

# Méthode MCMC

**Principe** : construire une chaîne de Markov qui admette  $\pi$  comme distribution limite

Un exemple de MCMC : l'algorithme de Metropolis-Hastings construit une chaîne de Markov  $\{\theta_k, k \in \mathbb{N}\}$

Initialisation  $\theta_0 \sim p$  et à chaque itération  $k = 1, 2, \dots$ , faire

- (i)  $\theta' \sim Q(\theta_{k-1}, \cdot)$
- (ii) définir  $\theta_k = \theta'$  avec proba.:

$$\begin{aligned}\alpha(\theta_{k-1}, \theta') &:= 1 \wedge \frac{\pi(\theta' | Y) Q(\theta_{k-1}, \theta')}{\pi(\theta_{k-1} | Y) Q(\theta_{k-1}, \theta')} \\ &= 1 \wedge \frac{f(Y_1, \dots, Y_N | \theta')}{f(Y_1, \dots, Y_N | \theta_{k-1})} \frac{p(\theta') Q(\theta_{k-1}, \theta')}{p(\theta_{k-1}) Q(\theta_{k-1}, \theta')}\end{aligned}$$

et  $\theta_k = \theta_{k-1}$  avec proba.  $1 - \alpha(\theta_{k-1}, \theta')$ .

**Problème** : la complexité de MH est en  $\mathcal{O}(N)$  par itération

## Exemple 2 (modèle déformable) : apprentissage fréquentiste

- ▶  $N = 10,000$  images  $Y_1, Y_2, \dots, Y_N$  et  $Y = \mathcal{M}_{15}\{(0,1)\}$
- ▶ pour chaque image  $Y_n$ , la classe  $I_n$  et le champ de déformation  $\beta_n$  sont inconnues
- ▶ le bruit additif est supposé Gaussien de sorte que

$$f(Y_1, \dots, Y_N, I_n = i, \beta_n = d | \theta) = \mathcal{N}(\Phi_{\beta_n} \theta_{i_n}, \sigma \text{Id}_{15})$$

- ▶ on pose une loi a priori sur les données manquantes  $(I_n, \beta_n) \sim g$
- ▶ la fonction à maximiser sur  $\Theta$  est

$$\theta \mapsto f(Y_1, \dots, Y_N | \theta) = \sum_i \int_d f(Y_1, \dots, Y_N, I_n = i, D_n = d | \theta) dg(i, d).$$

**Apprentissage fréquentiste** : l'algorithme Expectation-Maximization (EM) est une façon standard de trouver le maximum de  $f$  **en présence de données manquantes**

# Algorithme EM

**Principe** : construire une séquence de paramètres  $\{\theta_k, k \in \mathbb{N}\}$  qui converge vers un maximum de  $f$

Initialisation avec  $\theta_0 \in \Theta$  et à chaque itération  $k = 1, 2, \dots$  :

- (i) calculer l'espérance conditionnelle sous la loi des données manquantes sachant  $\theta_{k-1}$

$$S_k(\theta) = \mathbb{E} \{ \log f(Y_1, \dots, Y_N, X_1, \dots, X_N | \theta) | \theta_{k-1} \}$$
$$\stackrel{iid}{=} \sum_{n=1}^N \mathbb{E} \{ \log f(Y_n, X_n | \theta) | \theta_{k-1} \}$$

- (ii) définir  $\theta_k$  comme un maximiseur de  $\theta \mapsto S_k(\theta)$

**Problème** : la complexité de EM est en  $\mathcal{O}(N)$  par itération

# Plan

Quelques informations sur mon parcours

Statistique computationnelle : modélisation et apprentissage

**Apprentissage statistique en présence d'un grand nombre de données : stabilité et efficacité d'approximations de certains algorithmes.**

Quelques Illustrations

# Approximations d'algorithmes

**Thème de recherche actif** : conception d'algorithmes qui

- 1– "imitent" les algorithmes d'apprentissage traditionnels
- 2– ont une complexité plus faible que  $\mathcal{O}(N)$
- 3– héritent (dans une certaine mesure) des garanties théoriques de ces algorithmes

Nous présentons deux algorithmes qui approximent MH et EM

- ▶ description des algorithmes
- ▶ garanties théoriques
- ▶ outils mathématiques pour justifier l'intuition de ces algorithmes

## Approximation de MH (avec N. Friel, UCD et P. Alquier, ENSAE)

Simuler une chaîne de Markov  $\{\tilde{\theta}_k, k \in \mathbb{N}\}$  en utilisant un sous-ensemble des données  $Y_U \subset Y$  de taille  $n \ll N$  ( $U \subset \{1, \dots, N\}$ ) pour accepter/rejeter un nouvel état de la chaîne

$$\alpha(\theta, \theta') = 1 \wedge \frac{f(Y | \theta')}{f(Y | \theta)} \times \frac{p(\theta')Q(\theta', \theta)}{p(\theta)Q(\theta, \theta')}$$

$\xrightarrow{\text{remplacé par}} \hat{\alpha}(\tilde{\theta}, \theta' | U) = 1 \wedge \frac{f(Y_U | \theta')^{N/n}}{f(Y_U | \tilde{\theta})^{N/n}} \times \frac{p(\theta')Q(\theta', \tilde{\theta})}{p(\tilde{\theta})Q(\tilde{\theta}, \theta')}$

A chaque itération, le sous-ensemble  $Y_U$  est tiré suivant une loi

$$\nu_{n,\epsilon}(U) \propto \exp\{-\epsilon \|S(Y) - (N/n)S(Y_U)\|\}, \quad \epsilon > 0.$$

Pour un sous-ensemble  $Y_U \subseteq Y$ ,  $S(Y_U)$  est un vecteur de statistique résumé de  $Y_U$ .

Cette méthode approximant MH est appelée **ISS-MCMC**, acronyme pour *Informed Sub-Sampling MCMC*  
**Complexité** :  $\mathcal{O}(n)$

# Convergence de MH

On définit le niveau de convergence d'une chaîne de Markov par

$$D_k(\theta_0) := \|\Pr(\theta_k \in \cdot \mid \theta_0) - \pi\| .$$

MH produit une chaîne de Markov ergodique *i.e*

$$\forall \theta_0 \in \Theta, \quad \lim_{k \rightarrow \infty} D_k(\theta_0) = 0 \quad \text{p-s.}$$

et possiblement la convergence se fait géométriquement,

▶ simplement

$$\forall \theta_0 \in \Theta, \quad D_k(\theta_0) \leq C(\theta_0)\rho^k .$$

▶ uniformément

$$\sup_{\theta_0 \in \Theta} D_k(\theta_0) \leq C\rho^k .$$



# Héritage des propriétés de MH

Il n'y a aucune raison pour que ISS-MCMC soit ergodique :

$$\Pr(\tilde{\theta}_k \in A) \not\rightarrow_{k \rightarrow \infty} \pi(A)$$

## Hypothèse

*Il existe un vecteur de statistiques suffisantes  $S$  pour  $f$  telle que*

$$f(Y | \theta) = \psi(\theta) \exp\{\phi(\theta)^T S(Y)\}$$

## Proposition

1. *Nous avons*

$$\text{KL}(\pi(\cdot | Y), \pi(\cdot | Y_U)) \leq B(Y_U)$$

*et  $B(Y_U)$  est minimisé pour le(s) sous-ensemble(s)  $Y_U$  qui minimise(nt)  $\|S(Y) - S(Y_U)\|$ .*

2. *Asymptotiquement quand  $N \rightarrow \infty$ , nous avons*

$$\text{KL}(\pi(\cdot | Y), \pi(\cdot | Y_U))$$

*est minimisé pour le(s) sous-ensemble(s)  $Y_U$  qui minimise(nt)  $\|S(Y) - S(Y_U)\|$ .*

# Héritage des propriétés de MH pour ISS-MCMC

En l'absence de statistiques suffisantes, nous travaillons sous l'hypothèse suivante

## Hypothèse

Il existe un vecteur de statistique résumé  $S$  et  $\gamma > 0$  tels que

$$\forall \theta \in \Theta \quad \left\| \frac{N}{n} \log f(Y_U | \theta) - \log f(Y | \theta) \right\| \leq \gamma \left\| S(Y) - \frac{N}{n} S(Y_U) \right\|$$

## Proposition

- ▶ Si MH converge simplement géométriquement alors ISS-MCMC converge de la même façon vers une approximation de  $\pi$
- ▶ Si MH converge uniformément géométriquement alors ISS-MCMC converge de la même façon vers une approximation de  $\pi$  quantifiée par

$$\begin{aligned} & \|\pi - \tilde{\pi}_k\| \\ & \leq \kappa \sup_{\theta \in \Theta} \mathbb{E} \left\{ \frac{f(Y | \theta)}{f(Y_U | \theta)^{N/n}} \right\} \sup_{U \in \mathcal{U}_n} \mathbb{E} \left\{ \left| \frac{p(\tilde{\theta}) Q(\tilde{\theta}, \theta)}{p(\theta) Q(\theta, \tilde{\theta})} \left| \frac{f(Y_U | \theta)^{N/n}}{f(Y | \theta)} - \frac{f(Y_U | \tilde{\theta})^{N/n}}{f(Y | \tilde{\theta})} \right| \right| \right\} \end{aligned}$$

# Approximation de EM (avec É. Moulines, École Polytechnique)

## Hypothèse

Il existe un vecteur de statistique suffisante  $S$  i.e une fonction  $\Phi$  telle que

$$f(Y, X | \theta) = \Phi(\theta, S(Y, X)).$$

Dans notre méthode, les observations  $Y_1, Y_2, \dots$  sont traitées une par une

Initialisation à  $\tilde{\theta}_0 \in \Theta$ , la séquence  $\{\tilde{\theta}_k, k \in \mathbb{N}\}$  est définie par

1. approximation étape E : mise à jour des statistiques suffisantes

$$s_k = s_{k-1} + \varrho_k (S(Y_k, X_k) - s_{k-1}), \quad X_k \sim \pi(\cdot | Y_k, \tilde{\theta}_{k-1}),$$

avec  $\varrho_k = \varrho_0 k^{-\alpha}$ ,  $\alpha \in (0, 1)$

2. approximation étape M : mise à jour du paramètre

$$\tilde{\theta}_{k+1} = \bar{\theta}(s_{k+1}) := \arg \max_{\theta \in \Theta} \Phi(\theta, s_{k+1}).$$

Cette méthode approximant EM est appelée **MCoEM**, pour *Monte Carlo online EM*

**Complexité** :  $\mathcal{O}(1)$

# Héritage des propriétés théoriques de l'EM

Dans sa forme traditionnelle, la séquence de l'EM converge pour tout  $\theta_0 \in \Theta$  vers

$$\Theta^* := \{\theta \in \Theta, \nabla_{\theta'} f(Y | \theta')|_{\theta} = 0\}$$

## Proposition

*Sous certaines hypothèses techniques (sur la topologie de  $\Theta$ ), la séquence  $\{\tilde{\theta}_k, k \in \mathbb{N}\}$  produite par MCoEM converge vers*

$$\Theta^\dagger := \{\theta \in \Theta, \nabla_{\theta'} \text{KL}(\pi^*, f(\cdot | \theta'))|_{\theta'=\theta} = 0\}.$$

*presque-sûrement quand  $k \rightarrow \infty$ .  $\pi^*$  est le processus génératif des données qui peut ne pas appartenir à  $\mathcal{F}$ .*

# Héritage des propriétés théoriques de l'EM

## Monotonicité

- ▶ À chaque itération EM vérifie  $f(Y | \theta_{k+1}) \geq f(Y | \theta_k)$
- ▶ À chaque itération MCoEM vérifie

$$\text{KL}(\pi^*, f(\cdot | \tilde{\theta}_{k+1})) \leq \text{KL}(\pi^*, f(\cdot | \tilde{\theta}_k)) + o(\|\tilde{\theta}_{k+1} - \tilde{\theta}_k\|), \quad \text{p-s}$$

## Taux de convergence

- ▶ Conditionnellement à  $(\theta_0, Y)$ , EM est déterministe mais la convergence se fait en  $1/\sqrt{n}$
- ▶ Conditionnellement à  $(\tilde{\theta}_0, Y)$ , MCoEM est stochastique et la convergence se fait en  $1/\sqrt{\varrho_n} \propto 1/\sqrt{n^{-\alpha}}$ ,  $\alpha \in (0, 1)$ .
- ▶ On peut de plus prouver la normalité asymptotique pour MCoEM:

$$\varrho_n^{-1/2}(\tilde{\theta}_n - \theta^*) \Rightarrow \mathcal{N}(0, \Sigma),$$

où  $\Sigma$  est la matrice de variance covariance d'un problème de Lyapunov défini dépendant de  $\text{Var}\{S(Y, X) | Y, \theta\}$ .

## Outils mathématiques pour établir des résultats en présence d'approximations (1/2)

Pour une chaîne de Markov, on utilise la théorie de perturbation des opérateurs linéaires:

- ▶ caractérisation de la chaîne par l'opérateur de transition

$$K(\theta, A) := \Pr(\theta_{n+1} \in A \mid \theta_n = \theta)$$

- ▶ la loi de probabilité de la chaîne est simplement l'itération de  $K$

Comment une approximation locale se propage-t-elle asymptotiquement?

Exemple de résultat:

### Proposition (Mitrophanov (AAP, 2005))

*Soit une chaîne  $\{\tilde{\theta}_n, n \in \mathbb{N}\}$  ayant pour transition  $\tilde{K}$ . Si  $K$  est uniformément ergodique alors pour  $n$  suffisamment grand, nous avons :*

$$\|\Pr\{\theta_n \in \cdot\} - \Pr\{\tilde{\theta}_n \in \cdot\}\| \leq \kappa_n \|K - \tilde{K}\|$$

$\Rightarrow \tilde{K}$  joue le rôle d'une version bruitée de l'opérateur de référence (le noyau de transition de MH dans notre cas)

## Outils mathématiques pour établir des résultats en présence d'approximations (2/2)

L'EM peut également être réécrit en terme d'opérateur

$$\theta_{n+1} = K(\theta_n), \quad K(\theta) = \bar{\theta}(\mathbb{E}\{S(Y, X) | \theta\})$$

$K$  est un opérateur non linéaire

⇒ les points stationnaires de l'EM sont dans  $\ker h$ ,  $h := K - \text{Id}$

**Approximation stochastique** (Robbins et Monro, 1951) : trouver les zéros d'une fonction  $h$  ne pouvant être évaluée que de façon bruitée  $\hat{h}(\theta) \approx h(\theta)$

Sous certaines hypothèses sur l'approximation  $\hat{h}$ , la séquence  $\{\tilde{\theta}_n, n \in \mathbb{N}\}$

$$\tilde{\theta}_{n+1} = \tilde{\theta}_n + \varrho_n \hat{h}(\tilde{\theta}_n),$$

où  $\varrho_n \searrow 0$ ,  $\sum_n \varrho_n = \infty$ ,  $\sum_n \varrho_n^2 < \infty$  alors :

$$\tilde{\theta}_n \rightarrow \{\theta \in \Theta, h(\theta) = 0\} \text{ p-s.}$$

**MCoEM** :  $h \xrightarrow{\text{remplacé par}} \hat{h} = \bar{\theta}(S(Y_n, X_n)) - \text{Id}, \quad (Y_n, X_n) \sim \delta_{y_n} \otimes \pi(\cdot | \tilde{\theta}_n, Y_n)$

# Plan

Quelques informations sur mon parcours

Statistique computationnelle : modélisation et apprentissage

Apprentissage statistique en présence d'un grand nombre de données : stabilité et efficacité d'approximations de certains algorithmes.

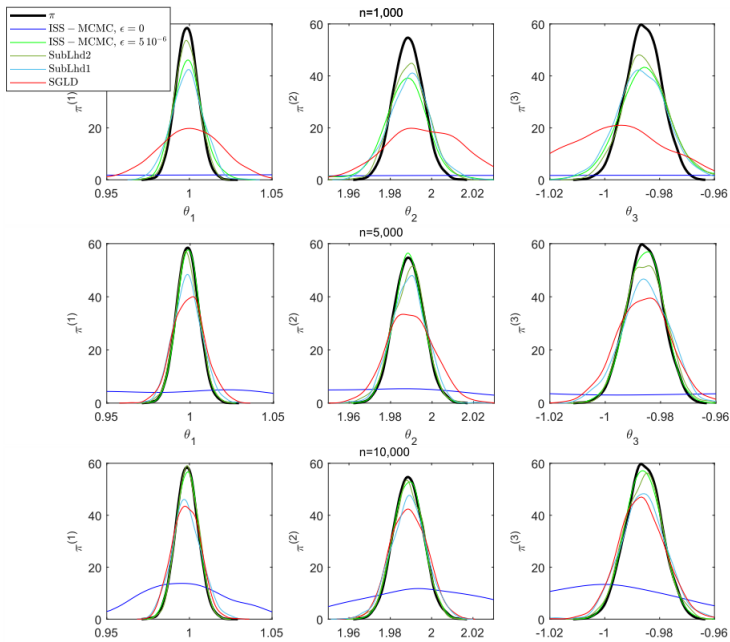
Quelques Illustrations



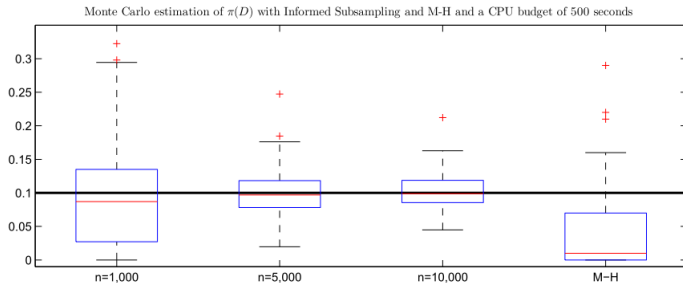
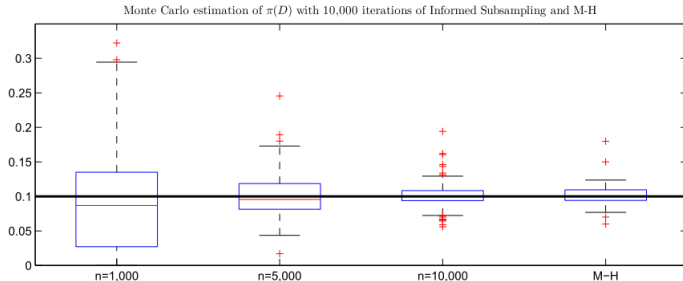
# Régression logistique (1/3)

algorithm	time/iter.(s)	iter. completed	RMSE	$\text{var}\{\widehat{\pi}(D)\}$
M-H	10	50	0.1417	0.004
ISS-MCMC, $n = 1,000$	0.05	10,000	0.1016	0.0104
ISS-MCMC, $n = 5,000$	0.08	6,250	0.0351	0.0012
ISS-MCMC, $n = 10,000$	0.13	3,840	0.0267	0.0007
SGLD, $n = 1,000$	0.08	6,000	0.1370	0.0157
SGLD, $n = 5,000$	0.11	5,250	0.0996	0.0100
SGLD, $n = 10,000$	0.12	4,500	0.0326	0.0011

# Régression logistique (2/3)



# Régression logistique (3/3)



# Modèle déformable

- Performance computationnelle

	MCoEM	EM (N=50)	EM (N=300)
temps CPU / itération (sec.)	20	225	1,570

- Performance en classification de chiffre manuscrit

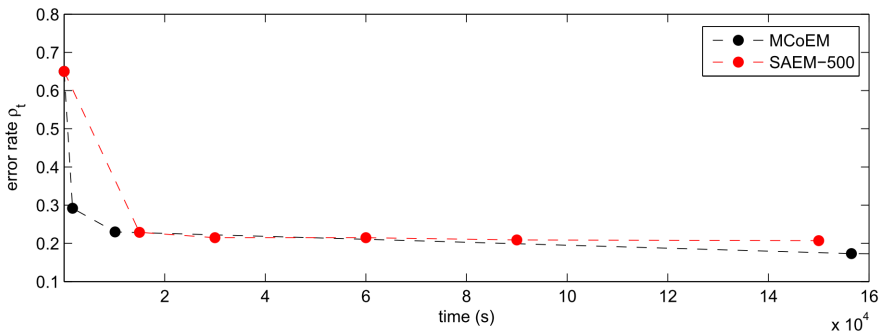


Figure: Erreur de classification en fonction du temps CPU

## Discussion

- ▶ Il est crucial d'établir des **approximations scalables** d'algorithmes bien connus afin d'exploiter au mieux les progrès des systèmes de mesure et de stockage
- ▶ Nous avons proposé de telles approximations pour **Metropolis-Hastings** et **Expectation-Maximization**
- ▶ Trouver des garanties théoriques pour ce genre de méthode représente un enjeu de taille, allant au delà de la communauté statistique
- ▶ Les résultats énoncés ici reposent sur **des hypothèses pas toujours réalistes** (existence de statistiques suffisantes ou "bonnes" statistiques résumées, convergence uniforme,  $\Theta$  compact...)
- ▶ D'autres outils théoriques ont été ou pourraient être utilisés comme le **couplage** de processus stochastiques ou la **théorie spectrale** des opérateurs