



Méthodes MCMC et applications Online EM Algorithm for Hidden Markov Model* Article d'Olivier Cappé

De Myttenaere Arnaud et Dadi Charles

6 janvier 2013

Introduction

L'algorithme EM a pour but initial de permettre d'estimer le maximum de vraisemblance dans un modèle dépendant de données latentes non observables. Ses applications sont très nombreuses dans le cadre notamment de l'imagerie médicale ou du traitement de la parole. La version classique de l'algorithme peut être qualifiée "d'hors ligne" dans la mesure où elle est exécutée à partir d'un échantillon d'observations statiques (dans le sens où cet échantillon n'évolue pas au cours de la procédure). Toutefois, l'apparition de problématiques liées au stockage des données de grande dimension, ou encore de l'adaptation de l'algorithme lorsque l'acquisition des données se fait de façon continue durant la procédure, conduit à modifier l'algorithme EM classique afin de prendre en compte les nouvelles observations au fur et à mesure de leur arrivée.

Dans l'article présenté ici [4], il s'agit de traiter des HMM (Hidden Markov Model), qui constituent une forme de modèle à données cachées, à savoir une chaîne de Markov de paramètres inconnus. On se restreint dans cet article à des HMM à espace d'états fini, avec une vraisemblance complète qui suit une distribution exponentielle.

*<http://arxiv.org/pdf/0908.2359.pdf>

Table des matières

1	De l’algorithme EM “hors” ligne à une version online	3
1.1	Rappel sur l’algorithme EM Classique :	3
1.2	Introduction de l’algorithme EM Online	3
1.3	Extension de l’algorithme EM Online à des données non indépendantes, HMM	3
1.4	Condition de convergence	4
1.5	Comparaison des performances	5
2	Implémentation et comparaison des algorithmes	6
2.1	EM classique pour un mélange de gaussiennes	6
2.2	EM online pour un mélange de gaussiennes	7
2.3	EM classique pour un HMM	8
2.4	EM online pour un HMM	11
3	Partie Théorique :	12
3.1	Vérification des hypothèses du théorème	12
3.2	Autour des preuves de fin d’article :	15
3.3	Utilisation des méthodes MCMC :	15
	Conclusion	16
	Références	16
4	Annexe	17
4.1	Code R pour l’algorithme EM Classique d’un mélange de gaussiennes	17
4.2	Algorithme EM online pour les HMM	19
4.3	Code R pour l’algorithme EM Classique de HMM	23

1 De l'algorithme EM "hors" ligne à une version online

1.1 Rappel sur l'algorithme EM Classique :

L'algorithme Maximisation-Expectation permet de maximiser une log-vraisemblance lorsque les observations dépendent d'une variable cachée. Ce modèle répond aux problématiques de données manquantes, d'observations bruitées.. L'algorithme se compose de deux étapes :

- l'étape (E) consiste à calculer l'espérance de la vraisemblance des données complètes.
- l'étape (M) permet ensuite de maximiser cette quantité afin d'actualiser les paramètres

L'idée générale pour justifier la convergence est la croissance de la log-vraisemblance à chaque itération. Cet algorithme se décline en de nombreuses versions adaptées à des modèles particulier tel que l'algorithme Forward-Backward pour estimer les paramètres d'une chaîne de Markov cachée (HMM - Hidden Markov Chain).

1.2 Introduction de l'algorithme EM Online

Dans le cadre de l'algorithme EM en ligne, les données sont considérées au fur et à mesure qu'elles sont disponibles. L'algorithme EM ne peut donc pas être lancé une fois toutes les données stockées. L'idée de base est de remplacer dans l'étape (E) le calcul de l'espérance de la vraisemblance complète par une étape d'approximation stochastique. La méthode, développée dans [4], propose une étape E différente. L'idée principale repose sur l'utilisation de la formule des probabilités totale pour approximer l'espérance, se ramenant ainsi au calcul de plusieurs espérances conditionnelles, en conditionnant par rapport au données observées. Comme dans de nombreux algorithmes de MCMC, l'itération nécessite maintenant l'introduction d'une suite $(\gamma_n)_{n \geq 0}$, telle que $\sum \gamma_n = +\infty$ et $\sum \gamma_n^2 < +\infty$. Après l'étape d'initialisation, les itérations s'effectuent le calcul au préalable des quantités suivantes :

$$\hat{\phi}_{n+1}(x) = \frac{\sum_{x' \in \mathcal{X}} \hat{\phi}(x') q_{\hat{\theta}_n}(x', x) g_{\hat{\theta}_n}(x, Y_{n+1})}{\sum_{x', x'' \in \mathcal{X}^2} \hat{\phi}(x') q_{\hat{\theta}_n}(x', x'') g_{\hat{\theta}_n}(x'', Y_{n+1})}$$

$$\hat{\rho}_{n+1}(x) = \sum_{x' \in \mathcal{X}} (\gamma_{n+1} s(x', x, Y_{n+1}) + (1 - \gamma_{n+1}) \hat{\rho}_n(x')) \frac{\hat{\phi}_n(x') q_{\hat{\theta}_n}(x', x)}{\sum_{x'' \in \mathcal{X}} \hat{\phi}_n(x'') q_{\hat{\theta}_n}(x'', x)}$$

L'espérance est alors estimée par :

$$\frac{1}{n} \widehat{\mathbb{E}} \left[\sum_{t=1}^n s(X_{t-1}, X_t, Y_t) | Y_0 : n \right] = \sum_{x \in \mathcal{X}} \hat{\phi}_{n+1}(x) \hat{\rho}_{n+1}(x)$$

Enfin, comme les données arrivent au fur et à mesure, il faut laisser passer un certain "temps de chauffe" avant d'actualiser l'estimation de θ pour la première fois. Dans la pratique, on fixe donc une valeur n_{min} , et pour $n > n_{min}$ pour effectuer l'étape M : $\hat{\theta}_{n+1} = \bar{\theta} \left(\sum_{x \in \mathcal{X}} \hat{\phi}_{n+1}(x) \hat{\rho}_{n+1}(x) \right)$

1.3 Extension de l'algorithme EM Online à des données non indépendantes, HMM

Dans le cadre de l'algorithme EM classique pour un modèle HMM, on suppose que (X_t, Y_t) sont générés suivant une chaîne de markov cachée de paramètre θ inconnu, avec l'hypothèse que X_t prend ses valeurs dans un espace fini \mathcal{X} et que seul les Y_i sont observés.

Une hypothèse importante de l'algorithme EM est le fait que le modèle est régi suivant une famille exponentielle, c'est-à-dire que la densité jointe des états et des observations s'écrit sous la forme :

$$p_\theta(x_t, y_t | x_{t-1}) = h(x_t, y_t) \cdot \exp(\langle \psi(\theta), s(x_{t-1}, x_t, y_t) \rangle - A(\theta))$$

L'étape M de l'algorithme EM consiste alors à déterminer le maximum de la vraisemblance pour un certain θ (évoluant à chaque étape et convergeant vers l'identification de la chaîne de Markov cachée), en résolvant pour $S \in \mathcal{S}$:

$$\nabla_\theta \psi(\theta) S - \nabla_\theta A(\theta) = 0$$

La solution est notée $\bar{\theta}(S)$.

En pratique, la $k^{ième}$ itération de l'algorithme EM, appliquée sur l'ensemble des observations, noté $Y_{0:n} = (Y_0, \dots, Y_n)$, est :

E - Step :

$$S_{k+1} = \frac{1}{n} \mathbb{E}_{\nu, \theta_k} \left[\sum_{t=1}^n s(X_{t-1}, X_t, Y_t) | Y_{0:n} \right]$$

M - Step : Mise à jour de l'estimation du paramètre θ : $\theta_{k+1} = \bar{\theta}(S_{k+1})$

L'étape E nécessite le calcul d'une espérance, ce qui n'est pas toujours possible en pratique. Il est donc envisageable d'utiliser des méthodes de MCMC, comme nous le verrons dans la partie. En 1988, Zeitouni et Dembo [9] proposent une méthode de calcul de cette espérance (en utilisant la formule des probabilités totales) à partir des quantités ϕ et ρ définies par :

$$\phi_{n, \nu, \theta}(x) = \mathbb{P}_{\nu, \theta}(X_n = x | Y_{0:n})$$

$$\rho_{n, \nu, \theta}(x) = \frac{1}{n} \mathbb{E} \left[\sum_{t=1}^n s(X_{t-1}, X_t, Y_t) | Y_{0:n}, X_n = x \right]$$

L'étape E peut alors être remplacée par le calcul de $\sum_{x \in \mathcal{X}} \phi_{n, \nu, \theta}(x) \rho_{n, \nu, \theta}(x)$. Cette méthode est aussi décrite dans [4]

1.4 Condition de convergence

Dans le cadre des HMM, la convergence est liée au théorème suivant, présenté dans [4] :

Si :

- i) la vraisemblance des données complète est de la forme exponentielle, et les paramètres sont des fonctions continument différentiable à l'intérieur de Θ ,
- ii) \mathcal{X} est un ensemble de cardinal fini,
- iii) Θ est un ensemble compact et θ^* est dans l'intérieur de Θ ,
- iv) la matrice de transition vérifie $q_\theta(x, x') \geq \varepsilon > 0$ pour tout $\theta \in \Theta$,
- v) $\sup_\theta \sup_y \bar{g}_\theta(y) < \infty$
- vi) $\mathbb{E}_{\theta^*} [\log \inf_\theta \bar{g}_\theta(Y_0)] < \infty$, où $\bar{g}_\theta(y) = \sum_x g_\theta(x, y)$

Alors :

i)

$$\frac{1}{n} \mathbb{E}_{\nu, \theta} \left[\sum_{t=1}^n s(X_{t-1}, X_t, Y_t) | Y_{0:n} \right] \xrightarrow{n \rightarrow \infty} \mathbb{E}_{\theta^*} [\mathbb{E}_{\theta} [s(X_{t-1}, X_t, Y_t) | Y_{-\infty; +\infty}]]$$

ii) Les points fixes de l'algorithme sont les points stationnaires de la vraisemblance.

La vérification de l'ensemble des hypothèses dans le cadre du modèle HMM présenté dans l'article fera l'objet de la partie 3 de ce rapport.

1.5 Comparaison des performances

Les expériences développées dans [4] montrent que sur les données utilisées l'algorithme EM en ligne converge, et que le choix de la suite γ_n influe sur la vitesse de convergence. Par ailleurs, la comparaison des résultats montre qu'à nombre d'observations identique l'algorithme classique propose des intervalles de confiance plus précis que l'algorithme en ligne, et est donc plus satisfaisant. Ce résultat se comprend bien dans la mesure où pour chaque itération l'algorithme EM classique utilise l'ensemble des observations, tandis que l'algorithme en ligne n'utilise que la nouvelle observation. Toutefois, les performances de l'algorithme EM online apparaissent très satisfaisantes.

Dans la partie suivante, nous reproduirons les résultats présentés dans [4] pour les HMM, et nous présenterons les résultats de l'estimation des paramètres dans le cas d'un mélange de gaussiennes pour les algorithmes EM classique et en ligne.

2 Implémentation et comparaison des algorithmes

2.1 EM classique pour un mélange de gaussiennes

Pour pouvoir souligner les avantages et les limites potentielles de l'algorithme EM Online, nous avons choisi de comparer les résultats obtenus avec un algorithme EM "hors" ligne. Dans un premier temps nous comparerons les deux algorithmes sur des données jouets échantillonnées selon un modèle de mélange de deux gaussiennes de même variance. Le modèle du mélange de deux gaussiennes répond à l'hypothèse d'indépendance des données (donc $L(y_{1:n}, x_{1:n}) = \prod_{t=1}^n L(y_t, x_t)$) et la vraisemblance complète se factorise sous la forme d'un modèle exponentiel de la façon suivante :

$$\begin{aligned}
 L(y_t, x_t) &= \mathbb{P}(y_t, x_t) = \mathbb{P}(y_t | x_t) \mathbb{P}(x_t) \\
 &= \prod_{i=1}^2 [\mathbb{P}(y_t | x_t = i) \mathbb{P}(x_t = i)]^{\mathbb{I}_{(x_t=i)}} \\
 &= \prod_{i=1}^2 \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{1}{2\sigma^2}(y_t - \mu_i)^2 + \log(p_i)\right]^{\mathbb{I}_{(x_t=i)}} \\
 &= \frac{1}{\sqrt{2\pi}} \exp\left[-\frac{1}{2} \log(\sigma) - \frac{1}{2\sigma} \sum_{i=1}^2 y_t^2 \mathbb{I}_{x_t=i} + \frac{1}{\sigma} \sum_{i=1}^2 y_t \mu_i - \frac{2}{\sigma} \mathbb{I}_{(x_t=i)} \sum_{i=1}^2 \mu_i \mathbb{I}_{(x_t=i)} + \sum_{i=1}^2 \log(p_i) \mathbb{I}_{(x_t=i)}\right]
 \end{aligned}$$

D'où,

$$s(y_{1:n}, x_{1:n}) = \left(\sum_{t=1}^n y_t^2 \mathbb{I}_{(x_t=1)}, \sum_{t=1}^n y_t^2 \mathbb{I}_{(x_t=2)}, \sum_{t=1}^n y_t \mathbb{I}_{(x_t=1)}, \sum_{t=1}^n y_t \mathbb{I}_{(x_t=2)}, \sum_{t=1}^n \mathbb{I}_{(x_t=1)}, \sum_{t=1}^n \mathbb{I}_{(x_t=2)} \right)$$

Ce modèle et l'algorithme ont été implémentés sous R (le code est mis en annexe). Les figures 1, 2 et 3 représentent respectivement les données et la convergence des variances et enfin la convergence des moyennes dans le cadre de l'algorithme EM "hors ligne" sur gaussiennes. Pour chaque graphique présenté tout au long de ce rapport, nous avons effectué 100 simulations, et nous avons représenté en rouge la valeur moyenne, et en vert premier et troisième quartiles, qui permettent de proposer des intervalles de confiance.

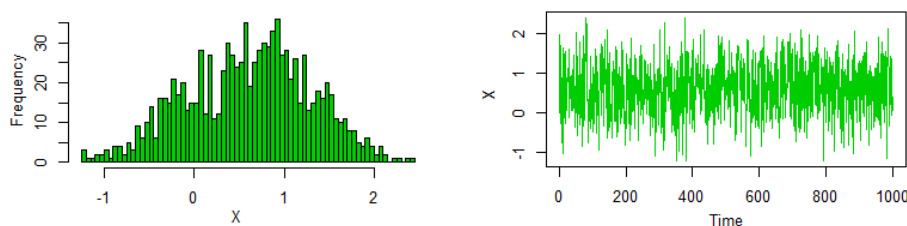


FIGURE 1 – Histogramme et trajectoire des observations générées selon un mélange de gaussiennes

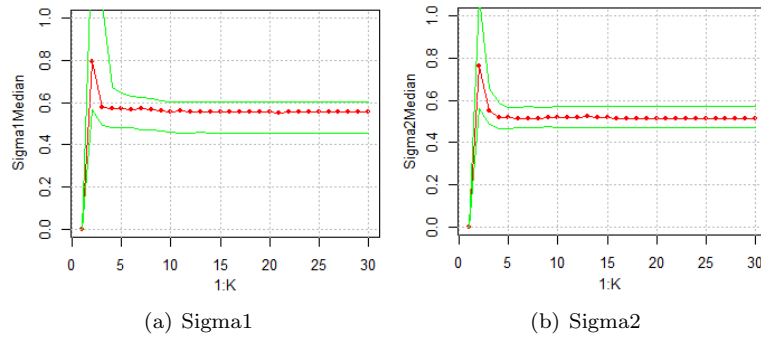


FIGURE 2 – Estimations de la variance des gaussiennes

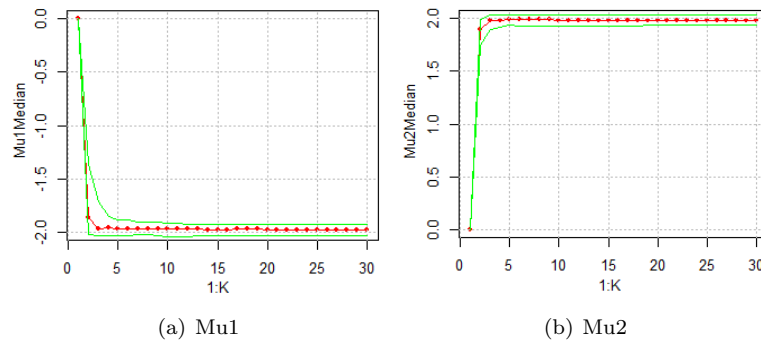


FIGURE 3 – Estimations de la moyenne des gaussiennes

On remarque une convergence très rapide (5 itérations) de l'algorithme EM classique. Pour confronter ces résultats à ceux de l'algorithme en ligne, intéressons-nous maintenant à l'implémentation online.

2.2 EM online pour un mélange de gaussiennes

Les figures 4 et 5 représentent la convergence de l'algorithme EM online pour le même mélange de gaussiennes que précédemment. Nous avons considéré une période de chauffe de 10 itérations, durant lesquelles les paramètres ne sont pas ajustés.

La convergence s'établit ici beaucoup plus lentement que dans le cas classique, puisque elle nécessite 30 à 40 itérations pour le cas en ligne. La période **constant** pour les premières itérations correspond à la période de temps de chauffe. Mais cette différence de vitesse de convergence s'explique par le fait que le nombre de données utilisées à chaque itérations est considérablement moindre dans le cas en ligne. En effet, l'estimation des paramètres ($\hat{\mu}$ et $\hat{\sigma}$) repose sur 20 observations à la 20^{ème} itération en ligne, alors qu'elle repose sur 100 observations à chaque itération de l'algorithme EM classique.

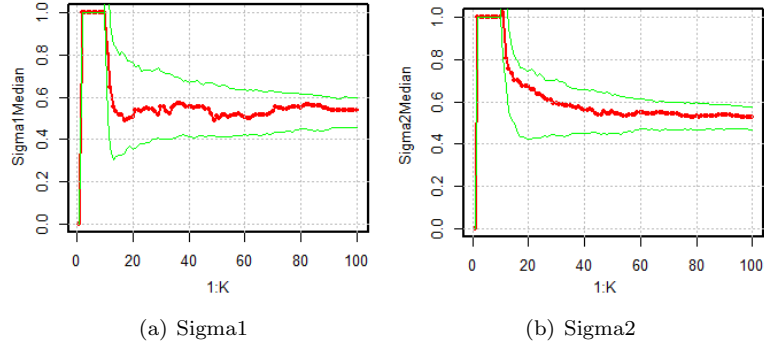


FIGURE 4 – Estimations de la variance des gaussiennes

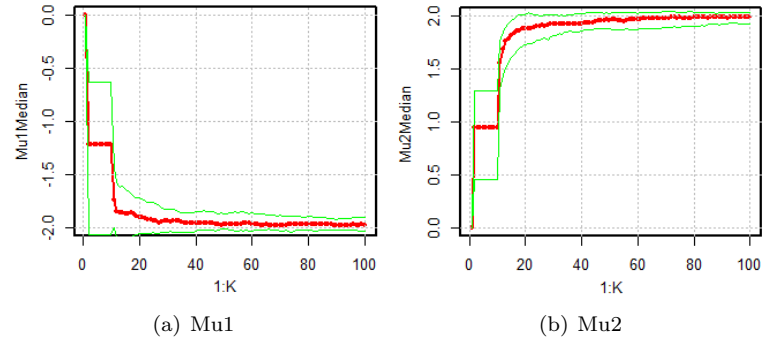


FIGURE 5 – Estimations de la moyenne des gaussiennes

Nous venons de voir que dans le cas des mélanges de gaussiennes les deux types d'algorithme EM convergent très rapidement (de 5 à 40 itérations). Cette vitesse de convergence, associée à la vitesse d'exécution des algorithmes, nous ont permis de réaliser plusieurs simulations pour obtenir des intervalles de confiance, et ainsi améliorer la précision de nos résultats. Etudions maintenant la convergence des algorithmes EM dans le cas d'un modèle à chaîne de Markov cachée.

2.3 EM classique pour un HMM

Le modèle EM et les algorithmes EM associés sont décrits dans [4]. Comme dans le modèle de mélange, le modèle HMM fait intervenir une structure sous jacente sous la forme d'une variable X_t . Une subtilité de l'article [4] est que la variable X_t est implicitement introduite, et qu'il s'agit d'une étiquette pointant vers une moyenne μ_i , ce qui permet en particulier d'alléger les notations.

On note ν la densité de X_0 et $\phi_k|n$ la densité de $X_k|Y_{0:n}$. On a alors

$$f_n(x, y, \theta) = \nu(x_0, \theta)g(x_0, x_1, \theta)q(x_0, x_1, \theta)g(x_1, y_1, \theta) \dots q(x_{n-1}, x_n, \theta)g(x_n, y_n, \theta)$$

L'expression de la log-vraisemblance est donc :

$$\log f_n(x_{0:n}, y_{0:n}, \theta) = \log \nu(x_0, \theta) + \sum_{t=0}^n \log(q(x_t, x_{t+1}, \theta)) + \sum_{t=0}^n \log g(x_t, y_{t+1}, \theta)$$

Nous nous intéressons à l'inférence dans ce modèle, et en particulier nous souhaitons calculer $\gamma_t(i) = \mathbb{P}(X_t = i | Y_{0:T})$, qui est la probabilité a posteriori de X_t sachant les données observées. Pour cela, nous avons utilisé la méthode introduite par Baum-Welch, appelée méthode *Forward Backward* [8], [6].

$$\begin{aligned}
\gamma_t(i) &= \mathbb{P}(X_t = i | Y_{0:T}) \\
&= \frac{\mathbb{P}(Y_{0:T} | X_t = i) \mathbb{P}(X_t = i)}{\mathbb{P}(Y_{0:T})} \\
&= \frac{\mathbb{P}(Y_{0:t} | X_t = i) \mathbb{P}(Y_{t+1:T} | X_t = i) \mathbb{P}(X_t = i)}{\mathbb{P}(Y_{0:T})} \\
&= \frac{\alpha_t(i) \beta_t(i)}{\mathbb{P}(Y_{0:T})}
\end{aligned}$$

On peut réécrire $\alpha_t(i)$, $\beta_t(i)$ et on obtient :

$$\begin{aligned}
\gamma_t(i) &= \frac{\alpha_t(i) \beta_t(i)}{\mathbb{P}(Y_{0:T})} \\
\alpha_{t+1}(i) &= \sum_{j=1}^2 f_i(y_{t+1}) \alpha_t(j) q_{i,j} \\
\beta_{t+1}(i) &= \sum_{j=1}^2 q_{i,j} f_j(y_{t+1}) \beta_{t+1}(j)
\end{aligned}$$

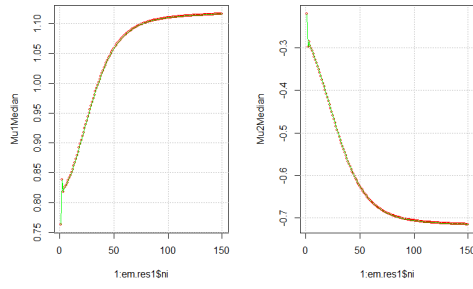
Ainsi nous pouvons calculer à présent les probabilités jointes pour la chaîne de Markov, $\mathbb{P}(X_t, X_{t+1} | Y)$ que l'on note dorénavant $\xi(x_t, x_{t+1})$. Grâce à la formule de Bayes et à la formule des probabilités totales on peut écrire :

$$\xi(i, j) = \frac{\alpha_t(i) \beta_{t+1}(j) q_{i,j} f_j(y_{t+1})}{\sum_{i=1}^2 \sum_{j=1}^2 \alpha_t(i) \beta_{t+1}(j) q_{i,j} f_j(y_{t+1})}$$

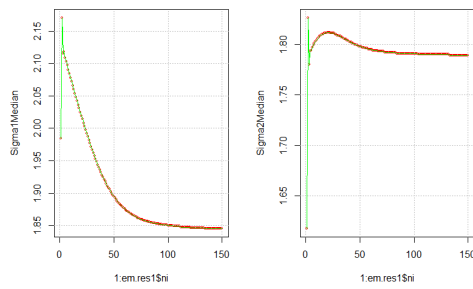
Remarque : en pratique, lors de l'implémentation nous avons été amenés à renormaliser les vecteurs α et β pour des problèmes d'*overflow*.

L'algorithme EM est alors le suivant :

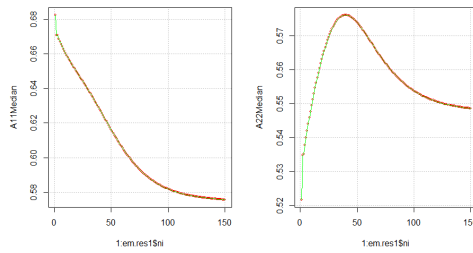
L'implémentation de l'algorithme EM "hors ligne" pour HMM se trouve en annexe. Les figures 6 représentent la convergence pour les termes de moyennes de variance et de la matrice de transition.



(a) Estimation des moyennes



(b) Estimation des variances



(c) Estimation de la diagonale de la matrice de transition

FIGURE 6 – Algorithme EM (version classique) appliqué à un modèle HMM

Un premier constat immédiat est la vitesse de convergence, qui est de l'ordre de 10 fois plus élevées que dans le cas du mélange de gaussienne. On peut donc s'attendre aussi à une vitesse de convergence beaucoup plus lente dans le cas de l'algorithme online.

2.4 EM online pour un HMM

Pour l'algorithme EM online dans le cas des HMM, nous avons d'abord généré une chaîne de Markov cachée à états, que nous avons ensuite bruitée pour obtenir des observations de la forme $Y_t = X_t + V_t$, où $V_t \sim \mathcal{N}(0, v^2)$. Puis nous avons estimé les valeurs de la matrice de transition q de la chaîne X_t , ainsi que ses états μ_1 et μ_2 . Enfin, comme précédemment, nous avons répété 100 fois cette expérience (de la génération des données à l'estimation des paramètres), afin de pouvoir tracer l'évolution des estimations moyennes et les premier et troisième quartiles, qui constituent en quelque sorte des intervalles de confiance.

Afin de comparer nos résultats avec ceux avancés dans l'article, nous avons utilisé les mêmes paramètres pour définir les observations, et les mêmes initialisations pour l'algorithme. Les observations sont donc définies par :

$$q_{1,1} = 0.95, q_{2,2} = 0.5, \mu_1 = 0, \mu_2 = 1 \text{ et } v = 0.5$$

Et l'algorithme est initialisé à :

$$q_{1,1}^{(0)} = 0.7, q_{2,2}^{(0)} = 0.5, \mu_1 = -0.5, \mu_2 = 0.5 \text{ et } v = 2$$

Les résultats pour la convergence de l'estimation de la matrice de transition et pour les états de la chaîne X_t sont présentés respectivement à la figure 7.

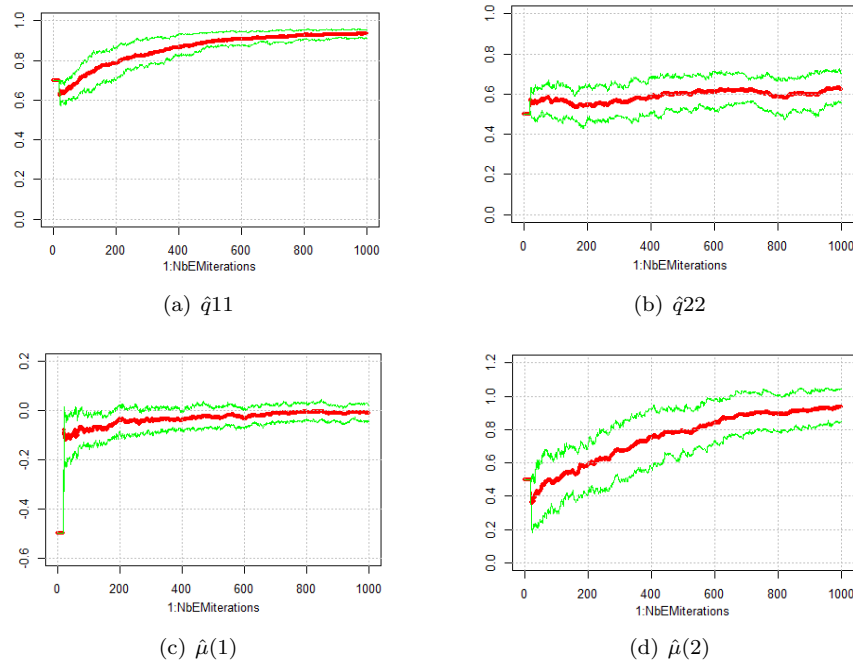


FIGURE 7 – Estimations des paramètres d'un HMM avec l'algorithme EM online

La vitesse de convergence dépend des états : pour q_{11} et μ_2 il faut attendre un grand nombre d'itérations (de l'ordre de 500, ce qui était prévisible d'après les résultats de l'algorithme EM classique), tandis que la convergence est presque acquise dès la sortie de la période de chauffe pour q_{22} et μ_1 .

3 Partie Théorique :

3.1 Vérification des hypothèses du théorème

Dans cette section nous vérifierons si le modèle *Markov Chain Observed in Gaussian noise* présenté en section 5.2 respecte les hypothèses du théorème 1, en page 8, énoncé dans [4]. Pour rappel, les théorèmes est énoncé dans la partie 1.4 de ce rapport.

Le modèle est le suivant : $Y_t = X_t + V_t$ où $V_t \rightsquigarrow N(0, v)$ et X_t est une chaîne de markov à deux états possibles $\mu(1), \mu(2)$. Pour l'expérimentation numérique nous avons posé les valeurs de θ initiales suivantes :

$$q_0 = \begin{pmatrix} 0.7 & 0.3 \\ 0.5 & 0.5 \end{pmatrix}, v_0 = 2, \mu_0(1) = -0.5 \text{ et } \mu_0(2) = 0.5$$

Hypothèse (i) : Vérifions que le modèle est exponentiel et la différentiabilité.

Calculons $\mathbb{P}[y_t, x_t | x_{t-1}]$:

L'attention du lecteur est attirée sur la variable x_t est la variable d'état, ou de régime. Elle prend ses valeurs dans $\{1, 2\}$ mais pointe sur les moyennes μ_1, μ_2 .

$$\begin{aligned} \mathbb{P}[y_t, x_t | x_{t-1}] &= \mathbb{P}[y_t | x_t, x_{t-1}] \mathbb{P}[x_t | x_{t-1}] \\ &= \prod_{i=1}^2 \prod_{j=1}^2 [\mathbb{P}(y_t | x_t, x_{t-1}) q_{i,j}]^{\mathbf{1}_{(x_t=i, x_{t-1}=j)}} \\ &= \frac{1}{\sqrt{2\pi}} \exp \left[\sum_{i=1}^2 \sum_{j=1}^2 \left(-\frac{1}{2} \log(\nu) - \frac{1}{2\nu} (y - \mu_i)^2 + \log(q_{i,j}) \right) \cdot \mathbf{1}_{(x_t=i, x_{t-1}=j)} \right] \\ &= \frac{1}{\sqrt{2\pi}} \exp \left[-2 \log(\nu) - \frac{1}{2\nu} \sum_{i=1}^2 2(y_t^2 - 2\mu_i y_t + \mu_i^2) \mathbf{1}_{(x_t=i)} + \sum_{i,j=1}^2 \log(q_{i,j}) \mathbf{1}_{(x_t=i, x_{t-1}=j)} \right] \end{aligned}$$

En notant $\Lambda = \nu^{-1}$, on a

$$\mathbb{P}[y_t, x_t | x_{t-1}] = \frac{1}{\sqrt{2\pi}} \exp \left[\log(\Lambda) - \Lambda \sum_{i=1}^2 y_t^2 \mathbf{1}_{(x_t=i)} + 2\Lambda \sum_{i=1}^2 \mu_i y_t \mathbf{1}_{(x_t=i)} - \Lambda \sum_{i=1}^2 \mu_i^2 \mathbf{1}_{(x_t=i)} + \sum_{i,j=1}^2 \log(q_{i,j}) \mathbf{1}_{(x_t=i, x_{t-1}=j)} \right]$$

Par identification, il s'agit d'un modèle exponentiel :

$$\mathbb{P}[y_t, x_t | x_{t-1}] = h(x_t, y_t) \exp [\psi(\theta), s(x_{t-1}, x_t, y_t) - A(\theta)]$$

$$\text{Où l'on a : } \psi(\theta) = \begin{pmatrix} \Lambda \\ \Lambda \\ \Lambda \mu(1) \\ \Lambda \mu(2) \\ \Lambda \mu(1)^2 \\ \Lambda \mu(2)^2 \\ \log(q_{1,1}) \\ \log(q_{1,2}) \\ \log(q_{2,1}) \\ \log(q_{2,2}) \end{pmatrix}, s(x_{t-1}, x_t, y_t) = \begin{pmatrix} -y^2 \mathbf{1}_{(x_t=1)} \\ -y^2 \mathbf{1}_{(x_t=2)} \\ 2y_t \mathbf{1}_{(x_t=1)} \\ 2y_t \mathbf{1}_{(x_t=2)} \\ \mathbf{1}_{(x_t=1)} \\ \mathbf{1}_{(x_t=2)} \\ \mathbf{1}_{(x_t=1, x_{t-1}=1)} \\ \mathbf{1}_{(x_t=1, x_{t-1}=2)} \\ \mathbf{1}_{(x_t=2, x_{t-1}=1)} \\ \mathbf{1}_{(x_t=2, x_{t-1}=2)} \end{pmatrix},$$

et $A(\theta) = 2 \log(\Lambda)$. Nous retrouvons ainsi les équations que respectent la statistique suffisante énoncées page 11 de [4].

Ainsi si on considère l'intérieur de Θ nous pouvons affirmer que ν^{-1} , $\mu_k \nu^{-1}$, $\log(q_{i,j})$ sont continument différentiables .

hypothèse (ii) : $\chi = \{0, 1\}$, donc χ est bien un ensemble de cardinal fini ($card(\chi) = 2$).

hypothèse (iii) :

$$\theta = \{q, \mu(1), \mu(2), \nu\} \in \mathcal{S}_{++}^2 \times \mathbb{R}^2 \times \mathbb{R}_+^*$$

Cet ensemble n'est pas compact. En effet, il suffit de considérer la suite $u_n = (q, \mu_1, \mu_2, \frac{1}{n})$.

$$\forall n \in \mathbb{N}^*, u_n \in \Theta, u_n \longrightarrow (q, \mu_1, \mu_2, 0) \notin \Theta$$

On peut cependant remarquer que cet ensemble est convexe. La notion de compacité nous assure une convergence de l'algorithme EM vers un point stationnaire de la log-vraisemblance. Nous savons que la vraisemblance est croissante avec le nombre d'itération de l'algorithme. Ainsi, supposons dans un premier temps l'existence de θ_* nous savons alors que $\forall i \in \{1, \dots, K\}$, $l(\theta_*) \geq l(\theta^i)$ où K est le nombre d'itération de l'algorithme EM.

Supposons à présent que $\theta^i \rightarrow \theta_*$ alors le fait que Θ soit compact nous assure que $\theta_* \in \Theta$. Donc pour toute fonction continue s , $\{s(\theta^i)\}_{i>0}$ converge vers $s(\theta_*)$. Mais ceci ne nous assure pas l'existence de θ_* ni la convergence de θ^i vers cette valeur. Toutefois la compacité de Θ nous assure par le théorème de Bolzano-Weierstrass qu'il existe au moins une sous suite de θ^i convergente.

Dorénavant on considérera un ensemble $\tilde{\Theta}$ compact inclus dans Θ . Si Θ est compact alors $\{\theta \in \Theta : l(\theta) \geq l(\tilde{\theta})\}, \forall \theta \in \Theta$.

hypothèse (iv) : A-t-on $q_\theta(x, x') \geq \varepsilon > 0 \forall \theta \in \Theta$? La réponse n'est pas évidente puisqu'en faisant tendre ν vers 0 on a un problème de définition. Il est cependant possible de montrer que la limite de $q(x, x')$ est strictement positive lorsque $\nu \rightarrow 0$, en remarquant que

$$\phi_{n+1}(k) = \frac{\sum_{k'} \phi_n(k') q_n(k', k) \psi_\theta(k, k^*, y, v)}{\sum_{k', k''} \phi_n(k') q_n(k', k'') \psi_\theta(k'', k^*, y, v)}$$

en notant k^* le k tel que $(y - \mu(k^*))^2 - (y - \mu(k))^2 \leq 0$ (possible car il y a seulement deux k), et

$$\psi(k, k', y, v) = \frac{g_\theta(k, y)}{g_\theta(k', y)} = \exp\left(\frac{(y - \mu(k'))^2 - (y - \mu(k))^2}{2v}\right)$$

Alors $\psi(k, k^*, y, 0) = 0$ pour $k \neq k^*$ et 1 sinon.

Donc

$$\lim_{\nu \rightarrow 0} \phi_{n+1}(k, \nu) = \begin{cases} 0 & \text{si } k \neq k^* \\ 1 & \text{sinon} \end{cases}$$

Et on peut montrer qu'à partir du rang 2 tous les ρ^q sont strictement positifs. Alors $S_{n+1}^q(i, j) = \sum_{k'} \rho_{n+1}^q(i, j, k') \phi_{n+1}(k') > 0$ pour tout (i, j) . Et donc $q_{n+1}(i, j) > 0$ pour tout (i, j) .

Si ν ne tend pas vers 0, remarquons que

$$\sum_{k=0}^m \phi(k) = 1 \text{ et } \forall k \in \{1, \dots, m\}, 0 \leq \phi(k) \leq 1$$

Donc il existe k tel que $\phi_n(k) \geq 1/m > 0$, et on en conclut comme précédemment que $q_{n+1}(i, j) > 0$ pour tout (i, j) .

En conclusion, nous n'avons pas $q_\theta(x, x') \geq \varepsilon > 0 \forall \theta \in \Theta$, mais nous pouvons affirmer que, sur tout compact inclut dans Θ , $\forall (x, x') \in \chi^2, q_\theta(x, x') \geq \varepsilon > 0$

hypothèse (v) : Montrons que $\sup_\theta \sup_y \bar{g}_\theta(y) < \infty$

On a :

$$\begin{aligned} \sup_\theta \sup_y \bar{g}_\theta(y) &= \sup_y \sup_\theta \bar{g}_\theta(y) \\ &= \sup_y \sup_\theta \sum_{x \in \{\mu(1), \mu(2)\}} g_\theta(x, y) \\ &\leq \sup_y \sum_{x \in \{\mu(1), \mu(2)\}} \sup_\theta g_\theta(x, y) \end{aligned}$$

Or ces densités étant gaussiennes nous avons

$$\forall \theta \in \Theta, g_\theta(\mu(1), y) + g_\theta(\mu(2), y) = \frac{1}{\sqrt{2\pi\nu}} e^{-\frac{1}{2\nu}(y-\mu(1))^2} + \frac{1}{\sqrt{2\pi\nu}} e^{-\frac{1}{2\nu}(y-\mu(2))^2}$$

Nous pouvons maximiser ces deux fonctions objectives par rapport à ν . On obtient une valeur maximale pour $\nu = (y - \mu(k))^2$.

Donc

$$\sup_y \sup_\theta \bar{g}_\theta(y) \leq \sup_y \left(\frac{1}{\sqrt{(2\pi(y - \mu(1))^2)}} + \frac{1}{\sqrt{(2\pi(y - \mu(2))^2)}} \right) \underset{p.s.}{<} \infty$$

hypothèse (vi) : A-t-on $\mathbb{E}_{\theta^*} [\log \inf_\theta \bar{g}_\theta(Y_0)] < \infty$?

Cette hypothèse n'est pas vérifiable car $\inf_\theta \bar{g}_\theta(Y_0) = 0$ (il suffit de faire tendre ν vers 0). En revanche, cette hypothèse est vraie pour tout compact inclut dans Θ . En effet :

Soit \mathcal{K} un compact inclut dans Θ .

La fonction $\bar{g}_\theta(Y_0)$ étant continue sur $\theta \in \mathcal{K}$ elle atteint ses bornes. Par ailleurs elle est strictement positive sur \mathcal{K} . Donc il existe c tel que $\bar{g}_\theta(Y_0) > c$, pour tout $\theta \in \mathcal{K}$. Alors

$$\inf_{\theta \in \mathcal{K}} \bar{g}_\theta(Y_0) > c$$

Par ailleurs, soit $\theta_0 \in \mathcal{K}$. On a :

$$\inf_{\theta \in \mathcal{K}} \bar{g}_\theta(Y_0) < \bar{g}_{\theta_0}(Y_0)$$

Ainsi

$$c < \inf_{\theta \in \mathcal{K}} \bar{g}_\theta(Y_0) < \bar{g}_{\theta_0}(Y_0)$$

Et $c > 0$, donc par croissance de la fonction \log :

$$\log(c) < \log\left(\inf_{\theta \in \mathcal{K}} \bar{g}_\theta(Y_0)\right) < \log(\bar{g}_{\theta_0}(Y_0))$$

Et donc

$$\mathbb{E} \left[\left| \log \left(\inf_{\theta \in \mathcal{K}} \bar{g}_\theta(Y_0) \right) \right| \right] \leq \max \left(\left| \log(c) \right|, \left| \log(\bar{g}_{\theta_0}(Y_0)) \right| \right)$$

D'où

$$\mathbb{E} \left[\left| \log \left(\inf_{\theta \in \mathcal{K}} \bar{g}_\theta(Y_0) \right) \right| \right] < \infty$$

Bilan : pour tout compact \mathcal{K} dans Θ , l'hypothèse est vraie.

3.2 Autour des preuves de fin d'article :

En annexe, [4] propose une preuve du théorème, par l'introduction de théorèmes énoncés dans [1], d'un lemme et d'un corollaire. Afin de comprendre la nécessité des hypothèses, nous nous sommes intéressés à la preuve du théorème. Plusieurs hypothèses apparaissent comme nécessaires dans les preuves :

- famille exponentielle et différentiabilité (théorème 3 de [1])
- le fait que χ soit de cardinal fini (corollaire 2),
- moments d'ordre 1 et 2 finis (théorème 1)
- θ^* est dans l'intérieur de Θ (théorème 1)

En revanche, la condition de compacité, qui est une hypothèse très forte et qui nous a posé plusieurs difficultés dans la partie vérification des hypothèses, n'apparaît pas clairement. Il est peut être possible d'alléger cette condition en supposant l'ensemble Θ convexe. Par ailleurs, les preuves nous apprennent que, dans un cadre plus général, une autre hypothèse peut être vérifiée par q et g simultanément : la condition présentée dans [1] et appelée *Strong Mixing Condition* :

Il existe un noyau de transition K et des fonctions mesurables ξ^- et ξ^+ telle que pour tout $(A, y) \in \xi \times \mathcal{Y}$,

$$\xi^-(y)K(y, A) \leq \int_A q(x, dx')g(x', y) \leq \xi^+(y)K(y, A)$$

3.3 Utilisation des méthodes MCMC :

Les méthodes de simulations MCMC permettent de générer des échantillons de la distribution a posteriori des paramètres du modèle. Dans l'étape E de l'algorithme EM Online il s'agit de calculer $\frac{1}{n} \mathbb{E}_{\nu, \theta_k} [\sum_{t=1}^n s(X_{t-1}, X_t, Y_t) | Y_{0:n}]$. On peut penser à utiliser une seule chaîne MCMC pour un nombre d'itérations pour laquelle la convergence est supposée atteinte et ainsi utiliser la dynamique de la chaîne pour échantillonner selon la loi a posteriori.

Cette méthode peut être une alternative intéressante à la simulation de plusieurs chaînes en parallèle, et améliorer sensiblement le temps de calcul.

Toutefois, il est nécessaire de poser des hypothèses sur le noyau de la chaîne. En effet, en simulant une seule chaîne on prend le risque de ne pas explorer tous les états. On peut par exemple supposer qu'une condition d'irréductibilité est nécessaire.

Conclusion

A la lumière de [4] et d'expérience portant sur les modèles de mélange de gaussiennes et de chaînes de Markov cachées, nous avons pu tester la pertinence de l'algorithme EM en ligne, comparée à l'algorithme EM classique. Il apparaît alors que la version en ligne est très satisfaisante et qu'elle permet une application de l'algorithme EM dans un contexte plus souple, puisqu'il n'est plus nécessaire de disposer de l'ensemble des données pour estimer un modèle. Il est cependant important de noter que la convergence est plus lente dans le cas de l'algorithme en ligne. La version en ligne nécessite donc moins d'observations mais un nombre d'itérations plus importants.

Toutefois, nous avons vu que la convergence de l'algorithme en ligne s'effectue par le théorème 1 sous de nombreuses hypothèses, parfois très restrictives (par exemple la compacité de Θ). L'analyse de ces hypothèses et les expériences montrent que ces conditions ne sont pas forcément nécessaires. Une réflexion nous a alors amené à penser que l'hypothèse portant sur la compacité peut peut-être être allégée, en supposant par exemple que l'ensemble des paramètres est convexe.

Au terme de cette étude il apparaît que l'algorithme EM Online permet une plus grande souplesse que l'algorithme EM classique, et son utilisation pourrait être préférée dans de nombreux domaines applicatifs pour faisant intervenir l'estimation de paramètres en temps réels, tel que le traitement automatique de la parole, ...

Références

- [1] Moulines et Ryden Cappé. Inference in hidden markov models. 1995.
- [2] Olivier Cappé. Version récursive de l'algorithme em pour l'estimation en ligne des paramètres de modèles de markov cachés.
- [3] Olivier Cappé. Modèles de mélange et modèles de markov cachés pour le traitement automatique de la parole. 2000.
- [4] Olivier Cappé. Online em algorithm for hidden markov chains. 2011.
- [5] Leggetter and Woodland. Maximum likelihood linear regression for speaker adaptation of continuous density hidden markov models. 1995.
- [6] MacKey. Ensemble learning for hidden markov models.
- [7] Eric Moulines Olivier Cappé and Tobias Ryden. Inference in hidden markov models. 2009.
- [8] Ganapathiraju Wu and Picone. Baum-welch re-estimation of hidden markov model. 1999.
- [9] Zeitouni and Dembo. Exact filters for the estimation of the number of transitions of finite-state continuous-time markov processes. 1988.

4 Annexe

4.1 Code R pour l'algorithme EM Classique d'un mélange de gaussiennes

```
1 #####
2 # Donnees issues dun melange gaussien
3 #####
4 GaussianSmoothing=function(N,p,mu1,mu2,sigma1,sigma2){
5   X=NULL
6   X<-matrix(data=0,nrow=N,ncol=1)
7   ### Boucle permettant de tirer N valeurs issues
8   ### dun melange gaussien :
9   for (i in 2:N) {
10     if (runif(1)<p) {X[i]=rnorm(1,mu1,sigma1)}
11     else {X[i]=rnorm(1,mu2,sigma2)}
12   }
13   X
14 }
15
16 X=GaussianSmoothing(N=1000,p=0.4,mu1=0,mu2=1,sigma1=0.5,sigma2=0.5)
17 par(mfrow=c(2,2))
18 hist(X,col="3",breaks=100)
19 ts.plot(X,col="3")
20
21 #####
22 # Algo EM : #
23 #####
24 #Definition et initialisation des parametres :
25 t1=proc.time()#horloge au debut de lalgo.
26 K = 20
27 #on fera K iterations de lalgorithme
28
29 hatP<-matrix(data=0,nrow=K,ncol=2) #matrice K,2 ..1ere colonne : p, 2nd:1-p
30 mu<-matrix(data=0,nrow=K,ncol=2)
31 sigma<-matrix(data=0,nrow=K,ncol=2)
32
33 hatP[1,1]=runif(1) #Initialisation de p uniforme
34 hatP[1,2]=runif(1)
35 mu[1,1] = runif(1,min=-1,max=1)
36 mu[1,2]= runif(1,min=1,max=2)
37 sigma[1,1] = runif(1,min=0,max=1)
38 sigma[1,2] = runif(1,min=0,max=1)
39
40 for (i in 2:K) {
41
42   vrais1 = hatP[i-1,1]*dnorm(X,mean=mu[i-1,1],sd=sigma[i-1,1])
43   vrais2 = (1-hatP[i-1,1])*dnorm(X,mean=mu[i-1,2],sd=sigma[i-1,2])
44
45   vrais12 = vrais1 / (vrais1 + vrais2) # probas a posteriori p-{i,1}
46   vrais22 = vrais2 / (vrais1 + vrais2) # probas a posteriori p-{i,2}
47   ## Mise a jour de lambda1 = P(Z=1 | X,Theta) :
48   hatP[i,1]=mean(vrais12) #p estime
49   hatP[i,2]= 1-hatP[i,1]
50
51   ## Mise a jour de mu1 et mu2 :
52   mu[i,1] = sum(vrais12*X)/sum(vrais12)
53   mu[i,2] = sum(vrais22*X)/sum(vrais22)
54   ## Mise a jour de sigma1 et sigma2 :
55   sigma[i,1] = sqrt(sum(vrais12*(X-mu[i,1])^2)/(sum(vrais12)))
```

```

57   sigma[i,2] = sqrt(sum(vrais22*(X-mu[i,2])^2)/(sum(vrais22)))
59 }
60 Mu1Median = rep(0,K)
61 Mu1Sup = rep(0,K)
62 Mu1Inf = rep(0,K)
63 for(i in 1:K){
64   s = summary(mu[,i])
65   Mu1Inf[i] = s[2]
66   Mu1Median[i] = s[4]
67   Mu1Sup[i] = s[5]
68 }
69
70 plot(1:K, sigma[,1], ylim=c(-2,2), pch = 16, type='o', col = "red", cex = 0.6)
71 plot(1:K, mu[,1], ylim=c(-2,2), pch = 16, type='o', col = "red", cex = 0.6)
72
73 t2=proc.time()-t1 # calcul le temps pour effecuter algo EM
74 t2
75 #boxplot(sigma[,1], main="Simple Box plot", col="lightblue")

```

4.2 Algorithme EM online pour les HMM

```

source("fonctions.R")
2
Nsimulations = 100;NbEMiterations = 1000
4
### Pour stocker les valeurs
6 Q11Array = array(0,dim = c(Nsimulations , NbEMiterations))

8 #####
### Simulations
10 #####

12 for(simulation in 1:Nsimulations){
  #####
  14 ### Generation des observations
  #####

  16 NbObs = 1000
  18 nmin = 20 ## Iteration a partir de laquelle on modifie les parametres

  20 Q<-matrix(data=c(0.95,0.3,0.05,0.7),nrow=2,ncol=2)# to define matrix Q
  mu=matrix(data=c(0,1),nrow=1,ncol=2) # to define vector mu

  22 Y=MarkovChain(N=NbObs,q=Q,mu1=mu[1],mu2=mu[2],sigma=.5)

  24 #####
  26 ### Algorithme
  #####

  28 Mu=c(-0.5,0.5)
  q = matrix(c(0.7,0.5,0.3,0.5), nrow = 2)
  30 v = 0.5
  nu = c(-0.5, 0.5) ### nu est la densite de X0 donc soit +- 0.5
  32 gamma = 1

  34 ### Initialisation des parametres
  phi = c( nu[1]*g(y,v)[1], nu[2]*g(y,v)[2] ) / (g(y,v)[1] + g(y,v)[2])
  36 Rho = array(0, dim = c(2,2,2))
  RhoD = array(0,dim = c(2,2,3))
  38 for(i in 1:2){
    for(k in 1:2){
      40 for(d in 1:3){
        RhoD[i,k,d] = delta(i-k)*(y^(d-1))
      }
    }
  }
  44 }

  46 Q11Array[simulation, 1] = q[1,1]

  48 ### C'est parti pour la boucle !
  for(n in 2:NbEMiterations){
    50 y=Y[n]
    gamma = n^(-0.6)

    52 ### Approx. Filter Update
    54 newPhi = actualiserPhi()

    56 ### E-Step
    newR = actualiserR()
    58 newRho = actualiserRho()
    newRhoD = actualiserRhoD()
  }

```

```

60     ### M-Step
61     if (n>=nmin){
62         newS = actualiserS()
63         q = actualiserQ()
64         newSD = actualiserSD()
65         Mu = actualiserMu()
66         v = actualiserV()
67     }
68
69     phi = newPhi
70     Rho = newRho
71     RhoD = newRhoD
72
73     Q11Array[simulation, n] = q[1,1]
74 }
75 }
76
77 Q11Median = rep(0, NbEMiterations)
78 Q11Sup = rep(0, NbEMiterations)
79 Q11Inf = rep(0, NbEMiterations)
80 for (i in 1:NbEMiterations){
81     s = summary(Q11Array[, i])
82     Q11Inf[i] = s[2]
83     Q11Median[i] = s[4]
84     Q11Sup[i] = s[5]
85 }
86
87 plot(1:NbEMiterations, Q11Median, main = "Q-11", ylim = c(0,1), pch = 16, type='o',
88      , col = "red", cex = 0.6)
89 grid(col = "grey", lty = "dotted", lwd = 1)
90 points(1:NbEMiterations, Q11Inf, pch = 20, col = "green", type='l', lwd = 1, cex =
91        1)
92 points(1:NbEMiterations, Q11Sup, pch = 20, col = "green", type='l', lwd = 1, cex =
93        1)

```

Les fonctions :

```

1 #####
2
3 ### Generation des donnees
4 #####
5 MarkovChain=function(N,q,mu1,mu2,sigma){
6     X=NULL
7     X<-matrix(nrow=N,ncol=1)
8     X[1,1]=mu2
9     V=rnorm(N,0,sigma)
10
11     for (i in 2:N){
12         u=runif(1)
13         if (X[i-1,1]==mu1){
14             if (u<q[1,1]){X[i,1]=mu1} else X[i,1]=mu2}
15         if (X[i-1,1]==mu2){
16             if (u<q[2,2]){X[i,1]=mu2} else X[i,1]=mu1}
17         }
18         Y=X+rnorm(N,0,sigma)
19         Y
20     }
21 }
22
23 g = function(y,v){

```

```

23 } c(exp(-(y - Mu[1])^2 / (2*v)), exp(-(y - Mu[2])^2 / (2*v)))
25 }
27 delta = function(x){
28   1*(x==0)
29 }
31 #####
32 ### Approx. Filter Update
33 #####
34 actualiserPhi = function(){
35   newPhi = c(0,0)
36   for(i in 1:2){
37     newPhi[i] = (phi[1]*q[1,i]*g(y,v)[i]+phi[2]*q[2,i]*g(y,v)[i])
38   }
39   newPhi/(sum(newPhi))
40 }
41 #####
42 ### E-Step
43 #####
44 actualiserR = function(){
45   newR = array(0, dim = c(2,2))
46   for(i in 1:2){
47     for(j in 1:2){
48       newR[i,j] = phi[i]*q[i,j]/(phi[1]*q[1,j] + phi[2]*q[2,j])
49     }
50   }
51   newR
52 }
53 }
54
55 actualiserRho = function(){
56   newRho = array(0, dim = c(2,2,2))
57   for(i in 1:2){
58     for(j in 1:2){
59       for(k in 1:2){
60         newRho[i,j,k] = gamma*delta(j-k)*newR[i,j] + (1-gamma)*(Rho[i,j,1]*newR[1,
61           k] + Rho[i,j,2]*newR[2,k])
62       }
63     }
64   }
65   newRho
66 }
67
68 actualiserRhoD = function(){
69   newRhoD = array(0, dim = c(2,2,3))
70   for(i in 1:2){
71     for(k in 1:2){
72       for(d in 1:3){
73         newRhoD[i,k,d] = gamma*delta(i-k)*(y^(d-1)) + (1-gamma)*(RhoD[i,1,d]*newR
74           [1,k] + RhoD[i,2,d]*newR[2,k])
75       }
76     }
77   }
78   newRhoD
79 }
80
81 #####
82 ### M-Step
83 #####
84 actualiserS = function(){

```

```

83 newS = array(0, dim = c(2,2))
84   for(i in 1:2){
85     for(j in 1:2){
86       newS[i,j] = newRho[i,j,1]*newPhi[1] + newRho[i,j,2]*newPhi[2]
87     }
88   }
89   newS
90 }
91
92 actualiserQ = function(){
93   q = array(0,dim = c(2,2))
94   for(i in 1:2){
95     for(j in 1:2){
96       q[i,j] = newS[i,j] / (newS[i,1] + newS[i,2])
97     }
98   }
99   q
100 }
101
102 actualiserSD = function(){
103   newSD = array(0, dim = c(2,3))
104   for(i in 1:2){
105     for(d in 1:3){
106       newSD[i,d] = newRhoD[i,1,d]*newPhi[1] + newRhoD[i,2,d]*newPhi[2]
107     }
108   }
109   newSD
110 }
111
112 actualiserMu = function(){
113   c(newSD[1,2] / newSD[1,1], newSD[2,2] / newSD[2,1])
114 }
115
116 actualiserV = function(){
117   numerateur = newSD[1,3] - (Mu[1]^2)*newSD[1,1] + newSD[2,3] - (Mu[2]^2)*newSD
118   [2,1]
119   numerateur / (newSD[1,1] + newSD[2,1])

```

4.3 Code R pour l'algorithme EM Classique de HMM

```
1  # Ce programme fonctionne a partir de 3 fonctions differentes
3  # On genere les donnees HMM a partir de 'rdata.hmm'
   #On realise la procedure foward backward grace a la fonction 'bwfw.hmm'
5  #Algorithme Em dans la fonction 'em.hmm'

7  # le nombre d observation
NUM1<-8000
9

11 #####
   #Parametres pour generer la chaine de Markov cachees
13 #####

15 # distribution initiale
   pii<-c(1, 0)
17 # matrice de transition
   A<-matrix(c(0.95, 0.05, 0.3, 0.70), 2, 2, byrow=T)
19 # 1ere distribution de y|x
   f0<-c(0, 2)
21 # 2nd distribution de y|x
   f1<-c(1, 2)
23 # simulation de HMM

25 rdata1<-rdata.hmm(NUM1, pii, A, f0, 1, f1)
   # donnees observees y
27 x1<-rdata1$o

29 par(mfrow=c(1,2))
   hist(rdata1$o,col="blue",breaks=100)
31 ts.plot(rdata1$o,col="blue")

33 # donnees latentes
   theta1<-rdata1$s
35

37 #####
   # Procedure forward/backward et scaling (etape E)
39 #####

41 x1<-rdata1$o
   fb.res1<-bwfw1.hmm(x1, pii, A, f0, f1)
43

45 # variable backward (beta)
   backward.var<-fb.res1$bw

47 # variable forward (alpha)
   forward.var<-fb.res1$fw
49

51 #####
   # Reestimation des paramatres (etape M)
53 #####
   L<-1

55 em.res1<-em1.hmm(x1, maxiter=150)

57 em.res1$A
   em.res1$f1
59 em.res1$pii
```

```

em.res1$muitier
61 em.res1$sditer
A11<-rep(0,em.res1$ni)
63 for(i in 1:em.res1$ni){A11[i]=em.res1$aiter[[i]][1,1]}
A22<-rep(0,em.res1$ni)
65 for(i in 1:em.res1$ni){A22[i]=em.res1$aiter[[i]][2,2]}

67 #On creer des tableau pour les graphiques de convergences.
Mu1Median = rep(0,em.res1$ni)
69 Mu1Sup = rep(0,em.res1$ni)
Mu1Inf = rep(0,em.res1$ni)
71 for(i in 1:em.res1$ni){
  s = summary(em.res1$muitier[i,1])
73   Mu1Inf[i] = s[2]
  Mu1Median [i] = s[3]
75   Mu1Sup [i] = s[5]
}
77 Mu2Median = rep(0,em.res1$ni)
Mu2Sup = rep(0,em.res1$ni)
79 Mu2Inf = rep(0,em.res1$ni)
for(i in 1:em.res1$ni){
81   s = summary(em.res1$muitier[i,2])
  Mu2Inf[i] = s[2]
83   Mu2Median [i] = s[3]
  Mu2Sup [i] = s[5]
85 }

87 Sigma1Median = rep(0,em.res1$ni)
Sigma1Sup = rep(0,em.res1$ni)
89 Sigma1Inf = rep(0,em.res1$ni)
for(i in 1:em.res1$ni){
91   s = summary(em.res1$sditer[i,1])
  Sigma1Inf[i] = s[2]
93   Sigma1Median [i] = s[3]
  Sigma1Sup [i] = s[5]
95 }

97
Sigma2Median = rep(0,em.res1$ni)
99 Sigma2Sup = rep(0,em.res1$ni)
Sigma2Inf = rep(0,em.res1$ni)
101 for(i in 1:em.res1$ni){
  s = summary(em.res1$sditer[i,2])
103   Sigma2Inf[i] = s[2]
  Sigma2Median [i] = s[3]
105   Sigma2Sup [i] = s[5]
}

107
A11Median = rep(0,em.res1$ni)
109 A11Sup = rep(0,em.res1$ni)
A11Inf = rep(0,em.res1$ni)
111 for(i in 1:em.res1$ni){
  s = summary(A11[i])
113   A11Inf[i] = s[2]
  A11Median[i] = s[3]
115   A11Sup [i] = s[5]
}

117
A22Median = rep(0,em.res1$ni)
119 A22Sup = rep(0,em.res1$ni)
A22Inf = rep(0,em.res1$ni)
121 for(i in 1:em.res1$ni){

```



```

123     s = summary(A22[i])
124     A22Inf[i] = s[2]
125     A22Median[i] = s[3]
126     A22Sup[i] = s[5]
127 }
128
129 #####
130 ### Graphiques
131 #####
132 par(mfrow=c(1,2))
133 #mul
134 plot(1:em.res1$ni, Mu1Median, type='o', col = "red", cex = 0.6)
135 grid(col = "grey", lty = "dotted", lwd = 1)
136 points(1:em.res1$ni, Mu1Inf, pch = 20, col = "green", type='l', lwd = 1, cex = 1)
137 points(1:em.res1$ni, Mu1Sup, pch = 20, col = "green", type='l', lwd = 1, cex = 1)
138
139 #mu2
140 plot(1:em.res1$ni, Mu2Median, type='o', col = "red", cex = 0.6)
141 grid(col = "grey", lty = "dotted", lwd = 1)
142 points(1:em.res1$ni, Mu2Inf, pch = 20, col = "green", type='l', lwd = 1, cex = 1)
143 points(1:em.res1$ni, Mu2Sup, pch = 20, col = "green", type='l', lwd = 1, cex = 1)
144
145 #Sigma1
146 plot(1:em.res1$ni, Sigma1Median, type='o', col = "red", cex = 0.6)
147 grid(col = "grey", lty = "dotted", lwd = 1)
148 points(1:em.res1$ni, Sigma1Inf, pch = 20, col = "green", type='l', lwd = 1, cex =
149 1)
150 points(1:em.res1$ni, Sigma1Sup, pch = 20, col = "green", type='l', lwd = 1, cex =
151 1)
152
153 #Sigma2
154 plot(1:em.res1$ni, Sigma2Median, type='o', col = "red", cex = 0.6)
155 grid(col = "grey", lty = "dotted", lwd = 1)
156 points(1:em.res1$ni, Sigma2Inf, pch = 20, col = "green", type='l', lwd = 1, cex =
157 1)
158 points(1:em.res1$ni, Sigma2Sup, pch = 20, col = "green", type='l', lwd = 1, cex =
159 1)
160
161 #A11
162 plot(1:em.res1$ni, A11Median, type='o', col = "red", cex = 0.6)
163 grid(col = "grey", lty = "dotted", lwd = 1)
164 points(1:em.res1$ni, A11Inf, pch = 20, col = "green", type='l', lwd = 1, cex = 1)
165 points(1:em.res1$ni, A11Sup, pch = 20, col = "green", type='l', lwd = 1, cex = 1)
166
167 #A22
168 plot(1:em.res1$ni, A22Median, type='o', col = "red", cex = 0.6)
169 grid(col = "grey", lty = "dotted", lwd = 1)
170 points(1:em.res1$ni, A22Inf, pch = 20, col = "green", type='l', lwd = 1, cex = 1)
171 points(1:em.res1$ni, A22Sup, pch = 20, col = "green", type='l', lwd = 1, cex = 1)
172
173 #####
174 #Fonction Backward/Forward et scaling
175 #####
176 bwfw1.hmm<-function(x, pii, A, f0, f1)
177 {
178
179 #####
180 ## USAGE

```

```

181 # bwfw.hmm(x, pii, A, f0, f1)
182
183 ## ARGUMENTS
184 # y=(y[1], ..., y[m]): donnee observees
185 # pii=(pii[0], pii[1]): distribution initiale
186 # A=(A[0,0], A[0,1]\ A[1,0], A[1,1]): matrice de transition
187 # f0=(mu, sigma): parametre pour y|x=1
188 # f1=(mu, sigma): parametre pour y|x=2
189
190 ## DETAILS
191 # bwfw.hmm calcule les valeurs forward et backward en tenant compte des pb d
192 # overflow
193 # --Procedure (Baum et al.)
194
195 #####
196
197 ## Initialisation
198
199 NUM<-length(y)
200
201 ## Densities
202
203 f0y<-dnorm(y, f0[1], f0[2])
204 f1y<-dnorm(y, f1[1], f1[2])
205
206 ## procedure Forward backward
207 alpha<-matrix(rep(0, NUM*2), NUM, 2, byrow=T)
208 # scaling variable c_0
209 c0<-rep(0, NUM)
210
211 alpha[1, 1]<-pii[1]*f0y[1]
212 alpha[1, 2]<-pii[2]*f1y[1]
213 #scaling
214 c0[1]<-1/sum(alpha[1, ])
215 alpha[1, ]<-c0[1]*alpha[1, ]
216
217 for (k in 1:(NUM-1))
218 {
219   alpha[k+1, 1]<-(alpha[k, 1]*A[1, 1]+alpha[k, 2]*A[2, 1])*f0y[k+1]
220   alpha[k+1, 2]<-(alpha[k, 1]*A[1, 2]+alpha[k, 2]*A[2, 2])*f1y[k+1]
221   # scaling
222   c0[k+1]<-1/sum(alpha[k+1, ])
223   alpha[k+1, ]<-c0[k+1]*alpha[k+1, ]
224 }
225
226
227 beta<-matrix(rep(0, NUM*2), NUM, 2, byrow=T)
228
229 beta[NUM, 1]<-c0[NUM]
230 beta[NUM, 2]<-c0[NUM]
231
232 for (k in (NUM-1):1)
233 {
234   beta[k, 1]<-A[1, 1]*f0y[k+1]*beta[k+1, 1]+A[1, 2]*f1y[k+1]*beta[k+1, 2]
235   beta[k, 2]<-A[2, 1]*f0y[k+1]*beta[k+1, 1]+A[2, 2]*f1y[k+1]*beta[k+1, 2]
236   # rescaling beta
237   # using the same scaling factors as alpha
238   beta[k, ]<-c0[k]*beta[k, ]
239 }
240
241 lfdR<-rep(0, NUM)

```

```

241 for (k in 1:NUM)
242 {
243   q1<-alpha[k, 1]*beta[k, 1]
244   q2<-alpha[k, 2]*beta[k, 2]
245   lfd r [k]<-q1/(q1+q2)
246 }
247
248 # probabilite des etats caches
249
250
251 gamma<-matrix(1:(NUM*2), NUM, 2, byrow=T)
252
253 gamma[NUM, ]<-c(lfd r [NUM], 1-lfd r [NUM])
254 dgamma<-array(rep(0, (NUM-1)*4), c(2, 2, (NUM-1)))
255
256 for (k in 1:(NUM-1))
257 {
258   denom<-0
259   for (i in 0:1)
260   {
261     for (j in 0:1)
262     {
263       fx<-(1-j)*f0x[k+1]+j*f1x[k+1]
264       denom<-denom+alpha[k, i+1]*A[i+1, j+1]*fx*beta[k+1, j+1]
265     }
266   }
267   for (i in 0:1)
268   {
269     gamma[k, i+1]<-0
270     for (j in 0:1)
271     {
272       fx<-(1-j)*f0x[k+1]+j*f1x[k+1]
273       dgamma[i+1, j+1, k]<-alpha[k, i+1]*A[i+1, j+1]*fx*beta[k+1, j+1]/denom
274       gamma[k, i+1]<-gamma[k, i+1]+dgamma[i+1, j+1, k]
275     }
276   }
277 }
278
279 # on retourne les resultats
280
281 bwfw.var<-list(bw=alpha, fw=beta, lsi=lfd r, pr=gamma, ts=dgamma)
282
283
284 return(bwfw.var)
285
286 }
287
288 #####
289 # ALGORITHME EM
290 #####
291 em1.hmm<-function(y, maxiter=50)
292 {
293
294 #
295 #####
296
297 ## ARGUMENTS
298 # y: donnees observee
299 # maxiter: Nombre d iteration maximum

```

```

301 # #####

303 #Tableau des stockages des valeurs des parametres au cours de EM
A_ITER<-vector("list",maxiter)
305 for(niter in 1:maxiter){A_ITER[[niter]] <- matrix(data = 0,nrow=2,ncol=2)}
MU_ITER<-matrix(data=0,nrow=maxiter,ncol=2)
307 SIGMA_ITER<-matrix(data=0,nrow=maxiter,ncol=2)
#####
309 NUM<-length(x)

311 niter<-0

313 ###initialisation des parametres
pii.new<-c(1, 0)
315 A.new<-matrix(c(0.7, 0.3, 0.5, 0.5), 2, 2, byrow=T)
f0.new<-c(-0.5, 2)
317 f1.new<-c(0.5, 2)

319

321 ### The E-M Algorithm

323 for(niter in 1:maxiter)
{
325

327 pii.old<-pii.new
A.old<-A.new
329 f1.old<-f1.new
f0.old<-f0.new
331
bwfw.res<-bwfw1.hmm(y, pii.old, A.old, f0.old, f1.old)
333
#probabilite des donnees cachees
335 gamma<-bwfw.res$pr
# the transition variables
337 dgamma<-bwfw.res$ts

339 ## (M STEP)

341 for (i in 0:1)
{
343 pii.new[i+1]<-gamma[1, i+1]
}
345

347 for (i in 0:1)
{
349 {
for (j in 0:1)
351 {
q1<-sum(dgamma[i+1, j+1, ])
353 q2<-sum(gamma[1:(NUM-1), i+1])
A.new[i+1, j+1]<-q1/q2
355 }
}
357 }

359 A_ITER[[niter]] = A.new

```

```

361 q11<-sum(gamma[, 2])
363 q22<-sum(gamma[, 2]*y)
    mu1<-q22/q11
365 q12<-sum(gamma[, 1])
    q22<-sum(gamma[, 1]*y)
367 mu2<-q22/q12
    MU_ITER[niter,]=c(mu1,mu2)
369 q31<-sum(gamma[, 2]*(y-mu1)*(y-mu1))
    q32<-sum(gamma[, 2]*(y-mu1)*(y-mu1))
371
373 sd1<-sqrt(q31/q11)
    sd2<-sqrt(q32/q12)
    SIGMA_ITER[niter,]=c(sd1,sd2)
375 f1.new<-c(mu1, sd1)
    f0.new<-c(mu2,sd2)
377 }
379 lfdr<-gamma[, 1]
381 em.var<-list(pii=pii.new, A=A.new, f1=f1.new, lf=lfdr, ni=niter, muiter=MU_ITER,
               sditer=SIGMA_ITER, aiter=A_ITER)
383 return(em.var)
385 }

```