
COMPTE-RENDU DU DEVOIR 3

Vincent Matthys

1 HMM implementation

1.0 Notations

On désignera par M le nombre d'états cachés (4 dans notre cas). On notera $\mathbf{u} = \{u_t\}_{t \in \llbracket 1; T \rrbracket}$, $\mathbf{q} = \{q_t\}_{t \in \llbracket 1; T \rrbracket}$, $\boldsymbol{\theta} = (\pi, A, \boldsymbol{\mu}, \boldsymbol{\Sigma})$ où $\boldsymbol{\mu} = \{\mu_k\}_{k \in \llbracket 1; M \rrbracket}$ et $\boldsymbol{\Sigma} = \{\Sigma_k\}_{k \in \llbracket 1; M \rrbracket}$. On utilisera également le *one-hot encoding* sur \mathbf{q} , de sorte que q_t^i , variable binaire, représente la i^{e} composante de q_t .

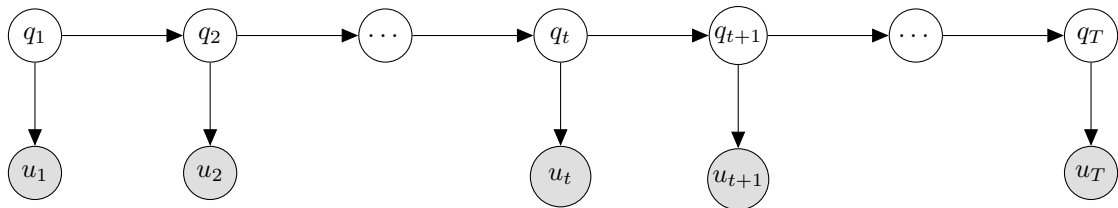


FIGURE 1 – Modèle graphique représentant le HMM considéré

1.1

Le problème de l'inférence dans la chaîne de Markov modélisée, à savoir le calcul de $p(\mathbf{q} \mid \mathbf{u})$ est appréhendable par le calcul de la marginale $p(q_t \mid \mathbf{u})$ suivant :

$$\begin{aligned}
p(q_t \mid \mathbf{u}) &= \frac{p(\mathbf{u} \mid q_t)p(q_t)}{p(\mathbf{u})} \\
&= \frac{p(u_1, \dots, u_t \mid q_t)p(u_{t+1}, \dots, u_T \mid q_t)p(q_t)}{p(\mathbf{u})} \quad u_1, \dots, u_t \perp\!\!\!\perp u_{t+1}, \dots, u_T \mid q_t \\
&= \frac{p(u_1, \dots, u_t, q_t)p(u_{t+1}, \dots, u_T \mid q_t)}{p(\mathbf{u})} \\
&= \frac{\alpha_t(q_t)\beta_t(q_t)}{p(\mathbf{u})} \\
&= \frac{\alpha_t(q_t)\beta_t(q_t)}{\sum_{q_t} \alpha_t(q_t)\beta_t(q_t)} \\
&= \frac{\alpha_t(q_t)\beta_t(q_t)}{\sum_{q_T} \alpha_T(q_T)} \quad \text{puisque} \quad p(\mathbf{u}) = \sum_{q_T} p(q_T, \mathbf{u}) = \sum_{q_T} \alpha_T(q_T)
\end{aligned} \tag{1}$$

avec $\alpha_t(q_t) = p(u_1, \dots, u_t, q_t)$ et $\beta_t(q_t) = p(u_{t+1}, \dots, u_T \mid q_t)$

On peut alors aussi exprimer $p(q_t, q_{t+1} \mid \mathbf{u})$ en fonction de α_t et β_t :

$$\begin{aligned}
p(q_t, q_{t+1} \mid \mathbf{u}) &= \frac{p(\mathbf{u} \mid q_t, q_{t+1})p(q_t, q_{t+1})}{p(\mathbf{u})} \\
&= \frac{p(\mathbf{u} \mid q_t, q_{t+1})p(q_{t+1} \mid q_t)p(q_t)}{p(\mathbf{u})} \\
&= \frac{p(u_1, \dots, u_t \mid q_t, q_{t+1})p(u_{t+1}, \dots, u_T \mid q_t, q_{t+1})p(q_{t+1} \mid q_t)p(q_t)}{p(\mathbf{u})} \\
&= \frac{p(u_1, \dots, u_t \mid q_t)p(u_{t+1}, \dots, u_T \mid q_{t+1})p(q_{t+1} \mid q_t)p(q_t)}{p(\mathbf{u})} \\
&= \frac{p(u_1, \dots, u_t \mid q_t)p(u_{t+1} \mid q_{t+1})p(u_{t+2}, \dots, u_T \mid q_{t+1})p(q_{t+1} \mid q_t)p(q_t)}{p(\mathbf{u})} \\
&= \frac{p(u_1, \dots, u_t, q_t)p(u_{t+1} \mid q_{t+1})p(u_{t+2}, \dots, u_T \mid q_{t+1})p(q_{t+1} \mid q_t)}{p(\mathbf{u})} \\
&= \frac{\alpha_t(q_t)p(u_{t+1} \mid q_{t+1})\beta_t(q_{t+1})A_{q_t, q_{t+1}}}{p(\mathbf{u})}
\end{aligned} \tag{2}$$

où $A_{q_t, q_{t+1}} = p(q_{t+1} \mid q_t)$ et $p(u_{t+1} \mid q_{t+1})$ suit une loi gaussienne.

À des fins d'implémentation, on considèrera les logarithmes des α et β , $\tilde{\alpha}$ et $\tilde{\beta}$ de sorte que les formules de récursion peuvent se réécrire :

$$\begin{cases} \ln \alpha_t(q_t) = \ln p(u_t \mid q_t) + \ln \sum_{q_{t-1}} (e^{\ln \alpha_{t-1}(q_{t-1})} A_{q_{t-1}, q_t}) \\ \ln \beta_t(q_t) = \ln \sum_{q_{t+1}} (e^{\ln \beta_{t+1}(q_{t+1})} A_{q_t, q_{t+1}} p(y_{t+1} \mid q_{t+1})) \end{cases} \tag{3}$$

$$\begin{cases} \tilde{\alpha}_t(q_t) = \ln p(u_t | q_t) + \ln \sum_{q_{t-1}} (e^{\tilde{\alpha}_{t-1}(q_{t-1})} A_{q_{t-1}, q_t}) \\ \tilde{\beta}_t(q_t) = \ln \sum_{q_{t+1}} (e^{\tilde{\beta}_{t+1}(q_{t+1})} A_{q_t, q_{t+1}} p(y_{t+1} | q_{t+1})) \end{cases} \quad (4)$$

et que la tâche d'inférence dite de *smoothing* en équation (1) s'écrit :

$$\begin{aligned} \ln p(q_t | \mathbf{u}) &= \ln \alpha_t(q_t) + \ln \beta_t(q_t) - \ln \sum_{q_T} \alpha_T(q_T) \\ \ln p(q_t | \mathbf{u}) &= \tilde{\alpha}_t(q_t) + \tilde{\beta}_t(q_t) - \ln \sum_{q_T} e^{\tilde{\alpha}_T(q_T)} \end{aligned} \quad (5)$$

La tâche d'inférence dite de *filtering* est également ajoutée en normalisant α_t :

$$\alpha_t^{\text{norm}}(q_t) = p(q_t | u_1, \dots, u_t) = \frac{\alpha_t(q_t)}{\sum_{q_t} \alpha_t(q_t)}$$

et est disponible par le calcul des α^{norm} dans le code fourni.

1.2

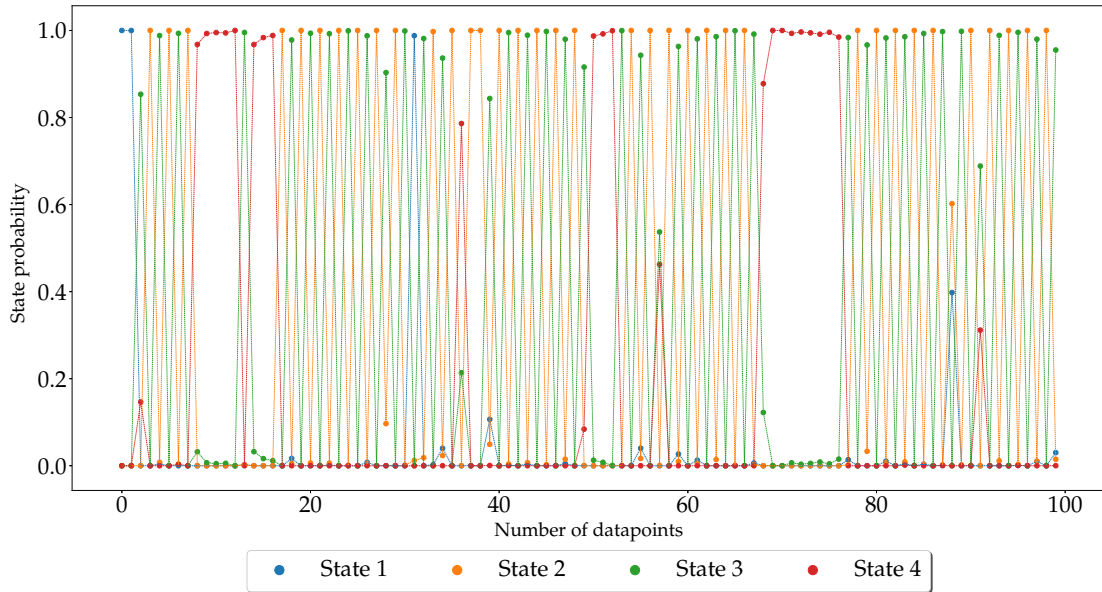


FIGURE 2 – Smoothing result for the first 100 datapoints

La figure donnant le résultat du smoothing des 100 premier points est présentée en figure 2. On y constate notamment une alternance importante entre les états 1 et 2, qui représentent également la majorité des états les plus probables.

1.3

On exprime la log-vraisemblance incomplète de la manière suivante :

$$\begin{aligned}
\ell(\boldsymbol{\theta}; \mathbf{u}) &= \ln \sum_{\mathbf{q}} p(\mathbf{q}, \mathbf{u} \mid \boldsymbol{\theta}) \\
&= \ln \sum_{\mathbf{q}} r(\mathbf{q} \mid \mathbf{u}) \frac{p(\mathbf{q}, \mathbf{u} \mid \boldsymbol{\theta})}{r(\mathbf{q} \mid \mathbf{u})} \quad \text{pour } r \text{ quelconque} \\
&\geq \sum_{\mathbf{q}} r(\mathbf{q} \mid \mathbf{u}) \ln \frac{p(\mathbf{q}, \mathbf{u} \mid \boldsymbol{\theta})}{r(\mathbf{q} \mid \mathbf{u})} \quad \text{par l'inégalité de Jensen} \\
&\geq \mathbb{E}_r[\ell_c(\boldsymbol{\theta}; \mathbf{q}, \mathbf{u})] - H(r(\mathbf{q} \mid \mathbf{u})) = \mathcal{L}(r, \boldsymbol{\theta}, \mathbf{u})
\end{aligned} \tag{6}$$

avec $\ell_c(\boldsymbol{\theta}; \mathbf{q}, \mathbf{u})$, log-vraisemblance complète, dont l'expression est la suivante :

$$\begin{aligned}
\ell_c(\boldsymbol{\theta}; \mathbf{q}, \mathbf{u}) &= \ln p(\mathbf{q}, \mathbf{u} \mid \boldsymbol{\theta}) \\
&= \sum_{t=1}^T \ln p(u_t \mid q_t, \boldsymbol{\mu}, \boldsymbol{\Sigma}) + \sum_{t=2}^T \sum_{i=1}^M \sum_{j=1}^M q_{t-1}^i q_t^j \ln A_{ij} + \sum_{i=1}^M q_1^i \ln \pi_i \\
&= \sum_{t=1}^T \sum_{k=1}^M q_t^k \left(-\ln 2\pi - \frac{1}{2} \ln |\boldsymbol{\Sigma}_k| - \frac{1}{2} (u_t - \mu_k)^\top \boldsymbol{\Sigma}_k^{-1} (u_t - \mu_k) \right) \\
&\quad + \sum_{t=2}^T \sum_{i=1}^M \sum_{j=1}^M q_{t-1}^i q_t^j \ln A_{ij} + \sum_{i=1}^M q_1^i \ln \pi_i
\end{aligned} \tag{7}$$

On a donc une borne inférieure de la log-vraisemblance incomplète, $\mathcal{L}(r, \boldsymbol{\theta}, \mathbf{u})$. L'algorithme EM consiste en une montée alternée de cette borne inférieure par rapport à r et $\boldsymbol{\theta}$. En choisissant $r(\mathbf{q} \mid \mathbf{u}) = p(\mathbf{q} \mid \mathbf{u}, \boldsymbol{\theta})$, on maximise celle-ci, de sorte que les deux étapes de l'algorithme se résument de la manière suivante :

$$\begin{aligned}
E - \text{step} \quad r^{(s+1)} &= \arg \max_r \mathcal{L}(r, \boldsymbol{\theta}^{(s)}, \mathbf{u}) = p(\mathbf{q} \mid \mathbf{u}, \boldsymbol{\theta}^{(s)}) \\
M - \text{step} \quad \boldsymbol{\theta}^{(s+1)} &= \arg \max_{\boldsymbol{\theta}} \mathcal{L}(r^{(s+1)}, \boldsymbol{\theta}, \mathbf{u}) = \arg \max_{\boldsymbol{\theta}} \mathbb{E}_{r^{(s+1)}}(\ell_c(\boldsymbol{\theta}; \mathbf{q}, \mathbf{u}))
\end{aligned} \tag{8}$$

où s dénote le nombre d'itérations de l'algorithme.

D'après l'équation (7), on observe que $m_{ij} := \sum_{t=2}^T q_{t-1}^i q_t^j$ est la statistique exhaustive pour A_{ij} , que q_1^i est la statistique exhaustive pour π_i . Les estimateurs de maximum de log-vraisemblance complète pour A et π sont donc :

$$\begin{aligned}
\hat{\pi}_i &= q_1^i \\
\hat{A}_{ij} &= \frac{m_{ij}}{\sum_{k=1}^M m_{ik}}
\end{aligned} \tag{9}$$

D'autre part, grâce aux calculs déjà effectués pour les estimateurs de vraisemblance de la mixture de gaussiennes, on obtient les estimateurs de log-vraisemblance complète des

paramètres d'émission :

$$\widehat{\mu}_k = \frac{\sum_{t=1}^T q_t^k u_t}{\sum_{t=1}^T q_t^k} \quad (11)$$

$$\widehat{\Sigma}_k = \frac{\sum_{t=1}^T q_t^k (u_t - \mu_k) (u_t - \mu_k)^\top}{\sum_{t=1}^T q_t^k} \quad (12)$$

Ces estimateurs doivent être pris sous l'espérance précise obtenue lors de l'étape E de l'algorithme EM, donnée par l'équation (8), afin que la maximisation de la log-vraisemblance complète induise une augmentation de la log-vraisemblance incomplète lors de l'itération $s + 1$.

En reprenant l'estimateur de π obtenu en équation (9), on obtient l'estimateur à l'itération $(s + 1)$:

$$\begin{aligned} \widehat{\pi}_i^{(s+1)} &= \mathbb{E}_{\mathbf{q}|\mathbf{u}, \boldsymbol{\theta}^{(s)}}[q_1^i] \\ &= \mathbb{E}_{q_1|\mathbf{u}, \boldsymbol{\theta}^{(s)}}[q_1^i] \\ &= p(q_1^i = 1 \mid \mathbf{u}, \boldsymbol{\theta}^{(s)}) := \gamma_1^{i, (s+1)} \end{aligned} \quad (13)$$

En reprenant l'estimateur de A_{ij} obtenue en équation (10), il faut déterminer $\mathbb{E}_{\mathbf{q}|\mathbf{u}, \boldsymbol{\theta}^{(s)}}[m_{ij}]$. On a :

$$\begin{aligned} \mathbb{E}_{\mathbf{q}|\mathbf{u}, \boldsymbol{\theta}^{(s)}}[m_{ij}] &= \sum_{t=2}^T \mathbb{E}_{\mathbf{q}|\mathbf{u}, \boldsymbol{\theta}^{(s)}}[q_{t-1}^i q_t^j] \\ &= \sum_{t=2}^T p(q_{t-1}^i = 1, q_t^j = 1 \mid \mathbf{u}, \boldsymbol{\theta}^{(s)}) \\ &:= \sum_{t=2}^T \xi_{t-1,t}^{ij, (s+1)} \end{aligned}$$

D'où l'estimateur à l'itération $(s + 1)$ pour A :

$$\widehat{A}_{ij}^{(s+1)} = \frac{\sum_{t=2}^T \xi_{t-1,t}^{ij, (s+1)}}{\sum_{k=1}^M \sum_{t=2}^T \xi_{t-1,t}^{ik, (s+1)}} = \frac{\sum_{t=2}^T \xi_{t-1,t}^{ij, (s+1)}}{\sum_{t=2}^T \gamma_{t-1}^{i, (s+1)}} \quad (14)$$

De la même façon, on obtient les estimateurs à l'itération $(s + 1)$ pour les μ_k et Σ_k :

$$\begin{aligned} \widehat{\mu}_k^{(s+1)} &= \frac{\sum_{t=1}^T \gamma_t^{k, (s+1)} u_t}{\sum_{t=1}^T \gamma_t^{k, (s+1)}} \quad \forall k \in \llbracket 1; M \rrbracket \\ \widehat{\Sigma}_k^{(s+1)} &= \frac{\sum_{t=1}^T \gamma_t^{k, (s+1)} (u_t - \mu_k^{(s)}) (u_t - \mu_k^{(s)})^\top}{\sum_{t=1}^T \gamma_t^{k, (s+1)}} \quad \forall k \in \llbracket 1; M \rrbracket \end{aligned} \quad (15)$$

L'implémentation du calcul des γ et ξ est succinctement détaillée et fait intervenir le

calcul des α^{norm} du problème de *filtering*.

$$\begin{aligned}
\gamma_t(q_t) &= p(q_t \mid u_1, \dots, u_T) \\
&= \sum_{q_{t+1}} p(q_t, q_{t+1} \mid u_1, \dots, u_T) \\
&= \sum_{q_{t+1}} \frac{\alpha_t(q_t) A_{q_t, q_{t+1}}}{\sum_{q_t} \alpha_t(q_t) A_{q_t, q_{t+1}}} \gamma_{t+1}(q_{t+1}) \\
&= \sum_{q_{t+1}} \frac{\alpha_t^{\text{norm}}(q_t) A_{q_t, q_{t+1}}}{\sum_{q_t} \alpha_t^{\text{norm}}(q_t) A_{q_t, q_{t+1}}} \gamma_{t+1}(q_{t+1})
\end{aligned} \tag{16}$$

$$\begin{aligned}
\xi_{t,t+1}(q_t, q_{t+1}) &= p(q_t, q_{t+1} \mid \mathbf{u}) \\
&= \frac{\alpha_t(q_t) p(u_{t+1} \mid q_{t-1}) \gamma_{t+1}(q_{t+1}) A_{q_t, q_{t+1}}}{\alpha_{t+1}(q_{t+1})} \\
&= \frac{\alpha_t^{\text{norm}}(q_t) p(u_{t+1} \mid q_{t-1}) \gamma_{t+1}(q_{t+1}) A_{q_t, q_{t+1}}}{\alpha_t^{\text{norm}}(q_{t+1})} \times \frac{\sum_{q_t} \alpha_t(q_t)}{\sum_{q_{t+1}} \alpha_{t+1}(q_{t+1})}
\end{aligned} \tag{17}$$

Au niveau de l'implémentation, on notera qu'on utilise les α^{norm} qui ne changent aucunement la formule pour γ . Pour ξ , il faut prendre en compte la différence de normalisation entre les α^{norm} .

1.4

La partie d'implémentation est effectuée toujours en python dans la classe *HMM*.

1.5

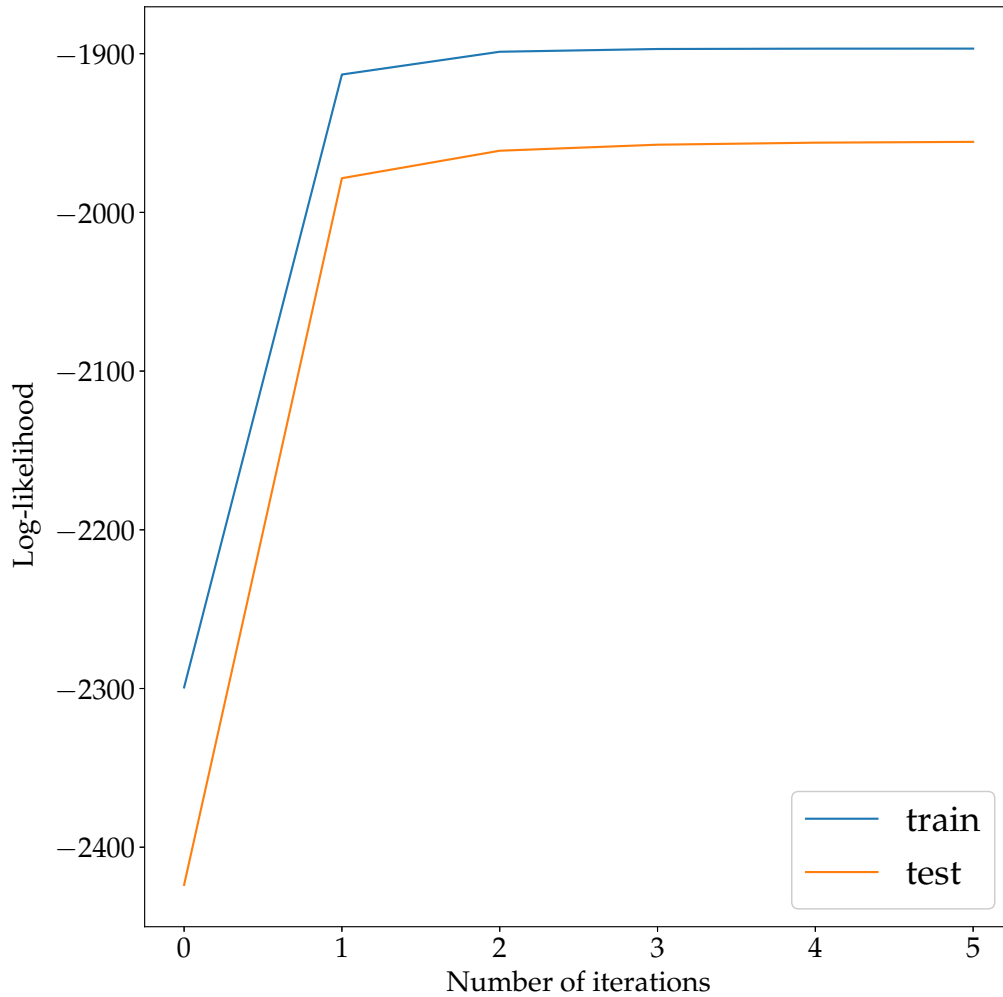


FIGURE 3 – Evolution de la log-vraisemblance en fonction du nombre d’itérations de EM pour le modèle de HMM

La log-vraisemblance en apprentissage et en test est représentée en figure 3. On constate qu’il suffit de 2-3 itérations de l’algorithme EM pour converger vers de nouveaux paramètres, en partant de ceux déterminés par mixture de gaussienne. On constate également que la vraisemblance en apprentissage est toujours supérieure à la vraisemblance de test, ce qui constitue un résultat attendu.

Modèle	Données d'entraînement	Données de test
isotropic GM	-2.6396	-2.6146
general GM	-2.3277	-2.4092
HMM	-1.8968	-1.9555

TABLE 1 – Valeur $\times 10^3$ de la log-vraisemblance des données d'apprentissage et des données de test pour différents modèles. GM signifie Gaussian Mixture (isotropic pour des covariances diagonales, general pour des covariances quelconques), HMM signifie Hidden Markov Model

1.6

Les valeurs de log-vraisemblance des différents modèles rencontrés sont présentées en table 1. On constate que le modèle avec la meilleure log-vraisemblance est le modèle chaîne de Markov, suivi du modèle de mixture de gaussiennes générales. C'est un résultat attendu, puisque l'ordre correspond avec le degré de complexité du modèle : plus celui-ci est complexe (c'est-à-dire avec de degré de liberté), plus il peut coller aux données. On constate également que la log-vraisemblance sur les données de test est inférieure à celle obtenue sur les données d'entraînement, à l'exception du modèle de mixture gaussiennes isotropiques pour lequel la différence observée est minime entre les données d'entraînement et de test. On notera cette vraisemblance de GM isotropic en test est la seule vraisemblance différente de celle notée dans la solution du devoir 2 fourni par M. Obozinski, dans laquelle la vraisemblance en question est de $-2.6841 \cdot 10^3$ et est bien inférieure à celle obtenue pour le même modèle sur les données d'entraînement.

1.7

L'algorithme de décodage de Viterbi permet de déterminer la séquence d'états la plus probable $\mathbf{q} = q_1, \dots, q_T$ étant donnée une séquence d'observation $\mathbf{u} = u_1, \dots, u_t$ pour des paramètres de modèle donnés $\boldsymbol{\theta}$. Au lieu de calculer la log-vraisemblance d'une séquence d'états étant donnée les observations, ce qui reviendrait à faire une α récursion pour chaque séquence d'états, on utilise l'algorithme *max-product* sur le graphe du modèle de Markov en figure 1, en prenant comme racine de l'arbre le noeud q_1 . Ainsi, on obtient la relation de récurrence suivante :

$$\delta_{t \rightarrow t+1}(q_{t+1}) = \max_{q_t} (\delta_{(t-1) \rightarrow t}(q_t) A_{q_t, q_{t+1}}) p(u_{t+1} | q_{t+1}) \quad (18)$$

$$\log \delta_{t \rightarrow t+1}(q_{t+1}) = \max_{q_t} (\log \delta_{(t-1) \rightarrow t}(q_t) + \log A_{q_t, q_{t+1}}) + \log p(u_{t+1} | q_{t+1}) \quad (19)$$

Une fois estimée ce message pour chaque point, on considère la succession d'états qui a permis d'obtenir ce message, notée $\mathbf{q}^* = \arg \max_{\mathbf{q}} p(\mathbf{q} | \mathbf{u})$, que l'on obtient en partant du noeud q_T :

$$q_T^* = \arg \max_{q_T} \delta_{(T-1) \rightarrow T}(q_T) \quad (20)$$

$$q_t^* = \arg \max_{q_t} \delta_{(t-1) \rightarrow t}(q_t) A_{q_t, q_{t+1}^*} \quad \forall t \in \llbracket 1; T-1 \rrbracket \quad (21)$$

D'où le pseudo-code suivant :

Algorithm 1: Viterbi decoding

Result: q^*

```

1  $\delta_1(j) = \pi_j p(u_1 \mid q_1 = j);$ 
2 for  $i = 2$  to  $T$  do
3    $\mid \delta_t(j) = \max_i (\delta_{t-1}(i) + \log A_{i,j}) + \log p(u_t \mid q_t = j);$ 
4 end
5  $q_T^* = \arg \max_j \delta_T(j);$ 
6 for  $i = T-1$  to  $1$  do
7    $\mid q_t^* = \arg \max_j \delta_t(j) A_{j, q_{t+1}^*};$ 
8 end
```

1.8

Le résultat du décodage de Viterbi sur les 100 premiers points des données d'entraînement avec les paramètres du modèle HMM obtenus par EM sur ces mêmes données est présenté en figure 4. On a également superposé les centres des clusters ainsi que l'ellipsoïde correspondant à 90 % de la masse de chaque gaussienne associée. On constate que l'attribution des points à chaque cluster est cohérent avec les résultats précédemment obtenus lors de la modélisation par mixture de gaussiennes générales.

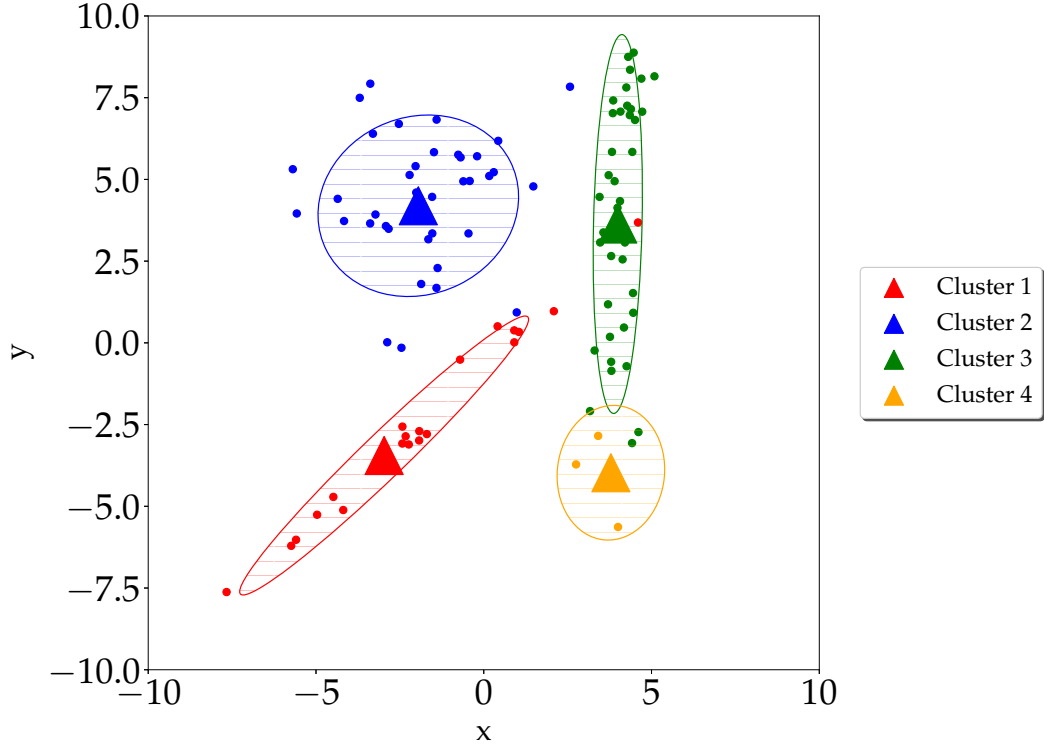


FIGURE 4 – Visualisation de la séquence d'états la plus probable sur les 100 premiers points du jeu d'entraînement obtenu par décodage de Viterbi

1.9

Le calcul des probabilités marginales de chaque état q_t étant donnée la séquence complète des observations \mathbf{u} a été décrite en équation (1) lors de l' $\alpha - \beta$ récursion. On utilise les paramètres α et β obtenus après convergence de EM sur les données d'entraînement, et on représente les marginales de chacun des 4 états des 100 premiers points du jeu de test en figure 5. On peut y observer une alternance presque exclusive des états 2 et 3.

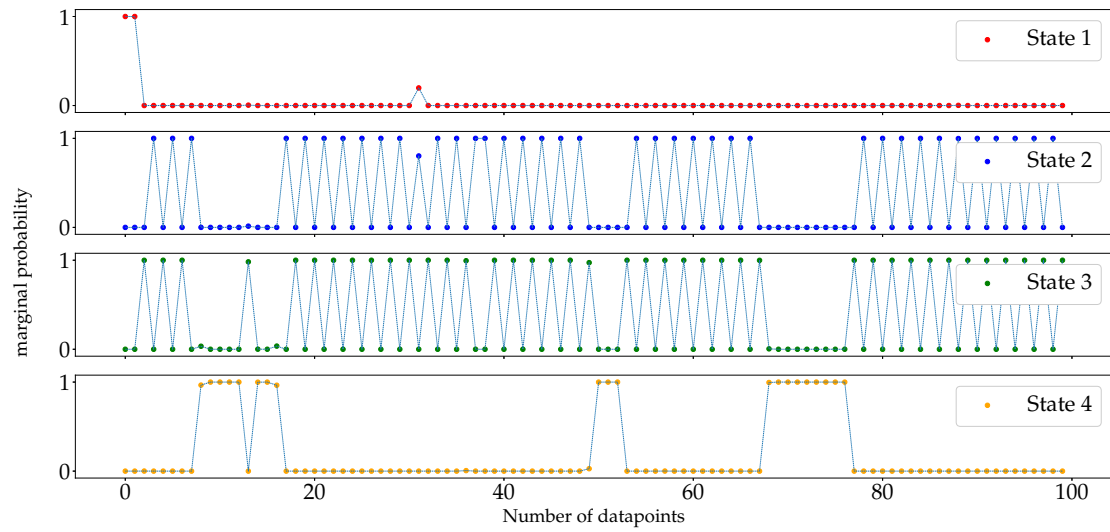


FIGURE 5 – Probabilités marginales des états sur le jeu de test pour les 100 premiers points

1.10

L'état le plus probable pour chacun de ces mêmes points est représenté en figure 6. On y constate clairement l'alternance entre les états 2 et 3, entrecoupés de périodes de 2 à 9 points en état 4.

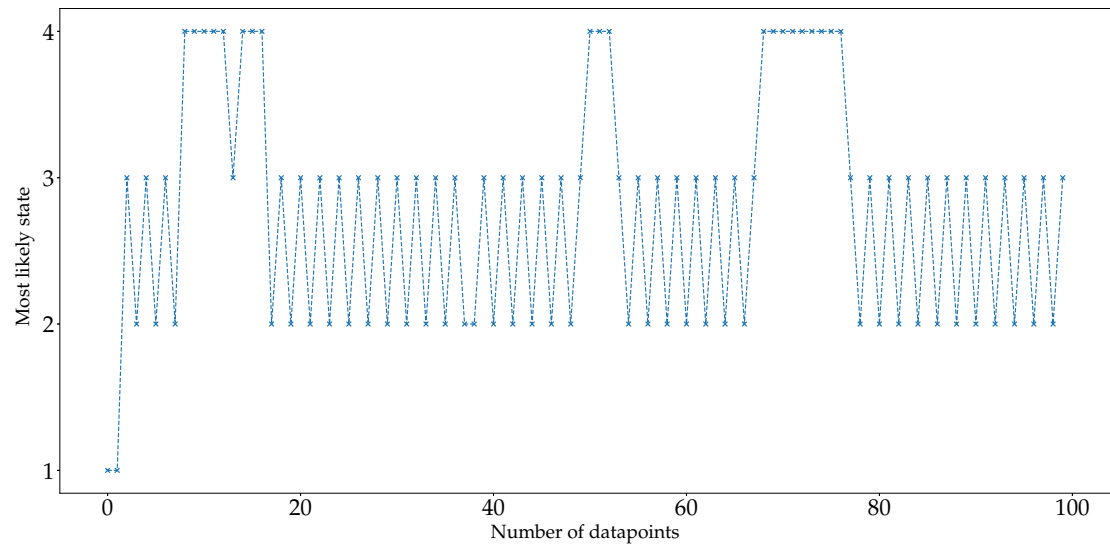
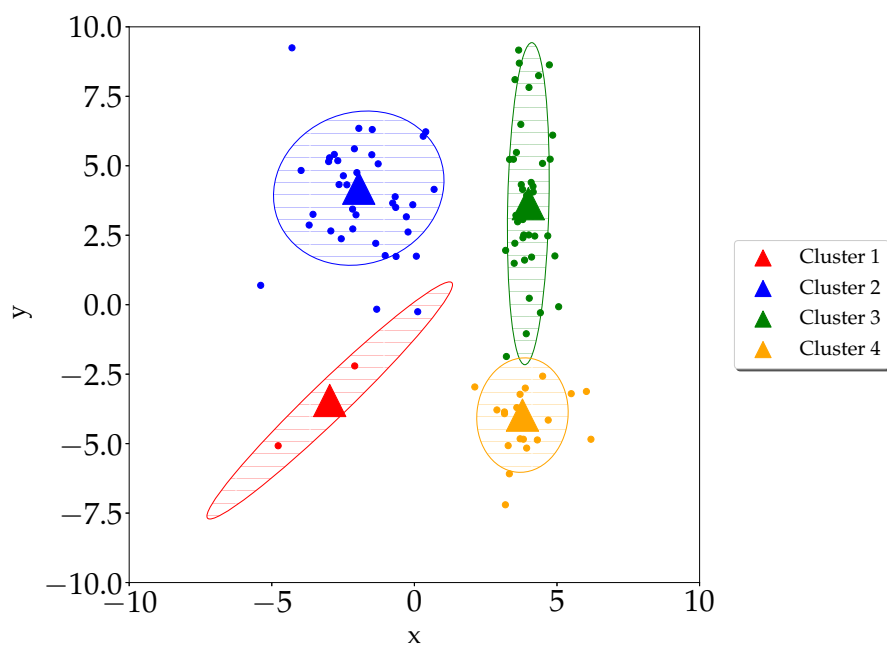
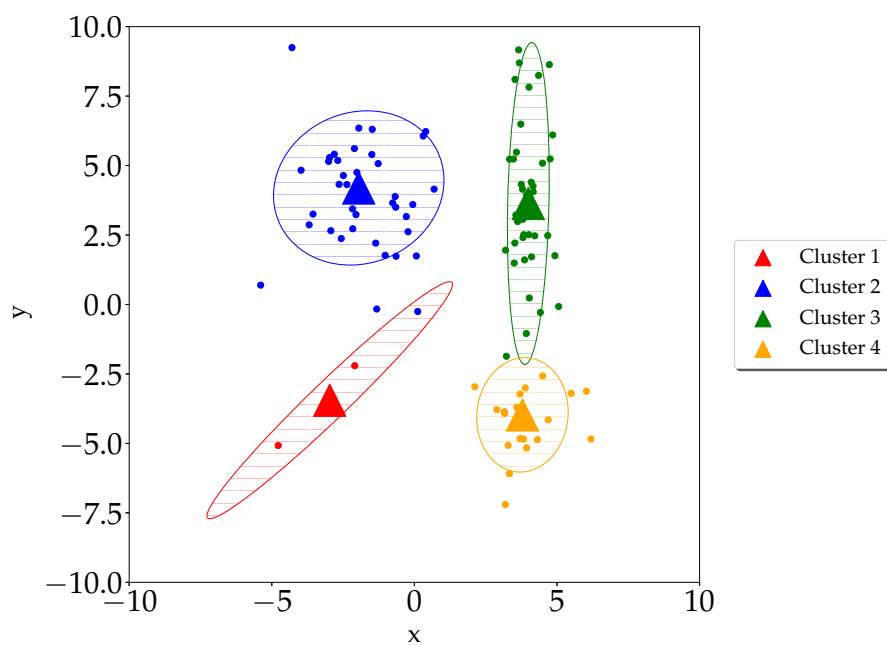


FIGURE 6 – Etat le plus probable pour chacun des 100 premiers points du jeu de test

1.11



(a) Décodage selon Viterbi



(b) Décodage selon les probabilités marginales

FIGURE 7 – Décodage des 100 premiers points du jeu de test suivant Viterbi ou les probabilités marginales

Les décodages obtenus par Viterbi et par les probabilités marginales pour les 100 premiers points du jeu de test sont présentés en figure 7. On remarque que les décodages par les deux méthodes sont identiques.

1.12

Si le nombre d'états n'était pas connu, il faudrait faire une sélection de modèle en comparant la log-vraisemblance obtenue par les différents modèles HMM (à différents nombre d'états) et conserver le modèle donnant la meilleure vraisemblance en test.