

L'OEP pour le problème de l'apprentissage de Modèles de Markov Cachés

Adaptation, hybridation et comparaison

Sébastien Aupetit, Nicolas Monmarché et Mohamed Slimane



Equipe Handicap et Nouvelles Technologies (HaNT)
Laboratoire d'Informatique
Université François Rabelais de Tours

Plan

- 1 Les Modèles de Markov Cachés (MMC)
- 2 L'apprentissage de MMC avec l'OEP
- 3 Etudes expérimentales
- 4 Conclusion

Les Modèles de Markov Cachés (MMC)

Les Modèles de Markov Cachés (MMC)

- Forme la plus simple de MMC : discret et en temps discret
- Outil mathématique permettant de modéliser une série temporelle par deux processus stochastiques :
 - un processus caché :
 - modélise les états (cachés) du système
 - c'est une chaîne de Markov
 - L'apparition d'un état ne dépend que de l'état précédent
 - un processus observé :
 - modélise l'émission de symboles observés
 - dépendant du processus caché
 - L'émission d'un symbole ne dépend que de l'état courant
- Paramètres :
 - Π : probabilités initiales des états cachés
 - A : probabilités de transition entre états cachés
 - B : probabilités d'émission des symboles par rapport aux états cachés

Les Modèles de Markov Cachés (MMC)

- Forme la plus simple de MMC : discret et en temps discret
- Outil mathématique permettant de modéliser une série temporelle par deux processus stochastiques :
 - un processus caché :
 - modélise les états (cachés) du système
 - c'est une chaîne de Markov
 - L'apparition d'un état ne dépend que de l'état précédent
 - un processus observé :
 - modélise l'émission de symboles observés
 - dépendant du processus caché
 - L'émission d'un symbole ne dépend que de l'état courant
- Paramètres :
 - Π : probabilités initiales des états cachés
 - A : probabilités de transition entre états cachés
 - B : probabilités d'émission des symboles par rapport aux états cachés

Les Modèles de Markov Cachés (MMC)

- Forme la plus simple de MMC : discret et en temps discret
- Outil mathématique permettant de modéliser une série temporelle par deux processus stochastiques :
 - un processus caché :
 - modélise les états (cachés) du système
 - c'est une chaîne de Markov
L'apparition d'un état ne dépend que de l'état précédent
 - un processus observé :
 - modélise l'émission de symboles observés
 - dépendant du processus caché
L'émission d'un symbole ne dépend que de l'état courant
- Paramètres :
 - Π : probabilités initiales des états cachés
 - A : probabilités de transition entre états cachés
 - B : probabilités d'émission des symboles par rapport aux états cachés

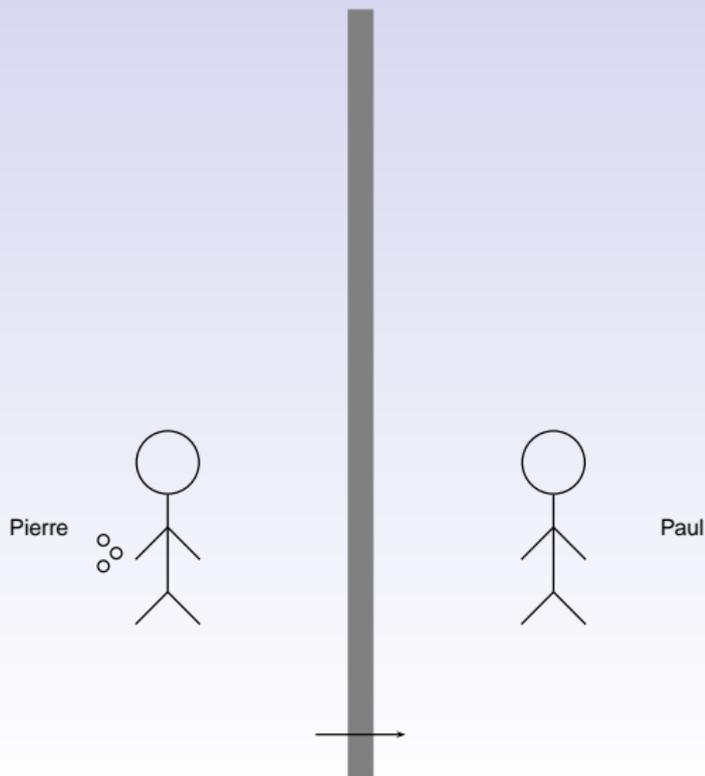
Les Modèles de Markov Cachés (MMC)

- Forme la plus simple de MMC : discret et en temps discret
- Outil mathématique permettant de modéliser une série temporelle par deux processus stochastiques :
 - un processus caché :
 - modélise les états (cachés) du système
 - c'est une chaîne de Markov
L'apparition d'un état ne dépend que de l'état précédent
 - un processus observé :
 - modélise l'émission de symboles observés
 - dépendant du processus caché
L'émission d'un symbole ne dépend que de l'état courant
- Paramètres :
 - Π : probabilités initiales des états cachés
 - A : probabilités de transition entre états cachés
 - B : probabilités d'émission des symboles par rapport aux états cachés

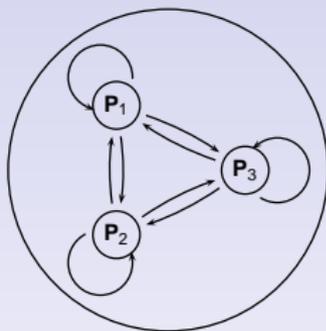
Les Modèles de Markov Cachés (MMC)

- Forme la plus simple de MMC : discret et en temps discret
- Outil mathématique permettant de modéliser une série temporelle par deux processus stochastiques :
 - un processus caché :
 - modélise les états (cachés) du système
 - c'est une chaîne de Markov
 - L'apparition d'un état ne dépend que de l'état précédent
 - un processus observé :
 - modélise l'émission de symboles observés
 - dépendant du processus caché
 - L'émission d'un symbole ne dépend que de l'état courant
- Paramètres :
 - Π : probabilités initiales des états cachés
 - A : probabilités de transition entre états cachés
 - B : probabilités d'émission des symboles par rapport aux états cachés

Exemple



Exemple



$$\Pi = (0.2 \ 0.4 \ 0.4)$$

$$A = \begin{pmatrix} 0.1 & 0.5 & 0.4 \\ 0.5 & 0.4 & 0.1 \\ 0.2 & 0.1 & 0.7 \end{pmatrix}$$

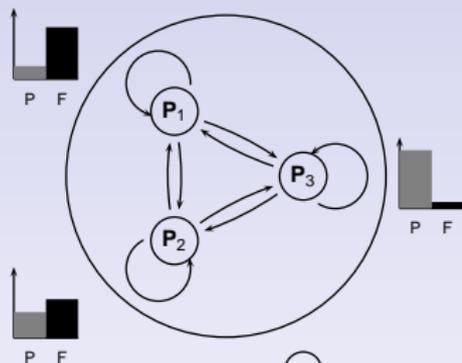
Pierre



Paul



Exemple

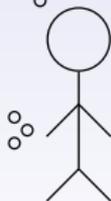


$$\Pi = (0.2 \ 0.4 \ 0.4)$$

$$A = \begin{pmatrix} 0.1 & 0.5 & 0.4 \\ 0.5 & 0.4 & 0.1 \\ 0.2 & 0.1 & 0.7 \end{pmatrix}$$

$$B = \begin{pmatrix} 0.2 & 0.8 \\ 0.4 & 0.6 \\ 0.9 & 0.1 \end{pmatrix}$$

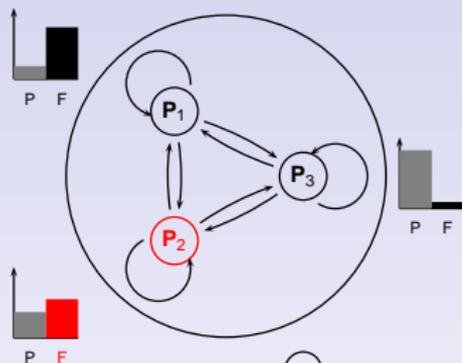
Pierre



Paul



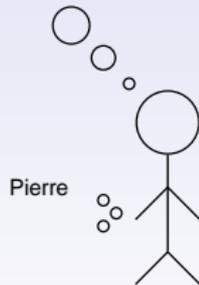
Exemple



$$\Pi = (0.2 \ 0.4 \ 0.4)$$

$$A = \begin{pmatrix} 0.1 & 0.5 & 0.4 \\ 0.5 & 0.4 & 0.1 \\ 0.2 & 0.1 & 0.7 \end{pmatrix}$$

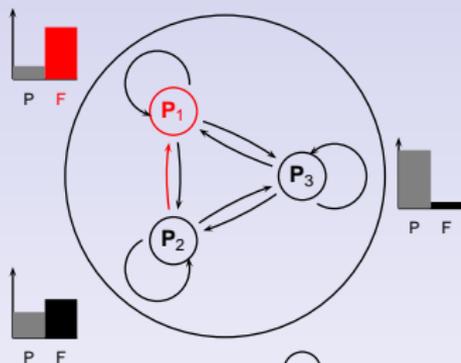
$$B = \begin{pmatrix} 0.2 & 0.8 \\ 0.4 & 0.6 \\ 0.9 & 0.1 \end{pmatrix}$$



P₂
F

F

Exemple

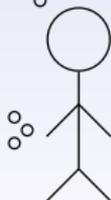


$$\Pi = (0.2 \ 0.4 \ 0.4)$$

$$A = \begin{pmatrix} 0.1 & 0.5 & 0.4 \\ 0.5 & 0.4 & 0.1 \\ 0.2 & 0.1 & 0.7 \end{pmatrix}$$

$$B = \begin{pmatrix} 0.2 & 0.8 \\ 0.4 & 0.6 \\ 0.9 & 0.1 \end{pmatrix}$$

Pierre



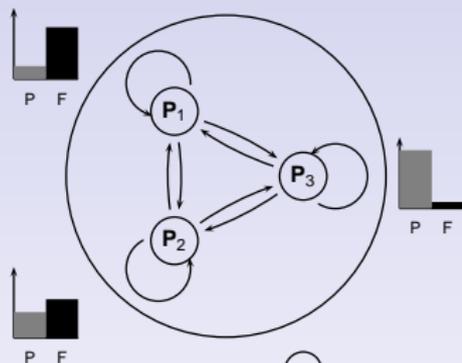
Paul



$P_2 P_1$
 $F F$

FF

Exemple

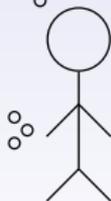


$$\Pi = (0.2 \ 0.4 \ 0.4)$$

$$A = \begin{pmatrix} 0.1 & 0.5 & 0.4 \\ 0.5 & 0.4 & 0.1 \\ 0.2 & 0.1 & 0.7 \end{pmatrix}$$

$$B = \begin{pmatrix} 0.2 & 0.8 \\ 0.4 & 0.6 \\ 0.9 & 0.1 \end{pmatrix}$$

Pierre



Paul



P₂P₁P₂P₂P₃P₁P₂P₁P₁P₃P₁P₂P₃P₁P₂P₃P₁
 F F P P P P F F P F F F P P P F



FFPPPPPPFFPPFFPPPPF

Les trois problèmes fondamentaux des MMC I

- Un modèle : $\lambda = (A, B, \Pi)$
- Trois problèmes fondamentaux :

- Calcul de la vraisemblance :

- $$P(O/\lambda) = \sum_{Q \in \mathbb{S}^T} P(O, Q/\lambda) = \sum_{Q \in \mathbb{S}^T} \left(\pi_{q_1} \left(\prod_{t=1}^{T-1} a_{q_t q_{t+1}} \right) \left(\prod_{t=1}^T b_{q_t}(o_t) \right) \right)$$

- Algorithmes Forward et Backward
- Calcul de la séquence d'états cachés la plus probablement suivie :
 - Trouver Q^* telle que $P(O, Q^*/\lambda)$ est maximale
 - Algorithme de Viterbi
- Apprendre une séquence d'observations :
 - Trouver λ^* qui maximise un critère
 - Nombreux critères : information mutuelle, discrimination, erreur de classification, ...

Les trois problèmes fondamentaux des MMC II

- Critère le plus utilisé :
 - Maximisation de la vraisemblance
 - Pas de solution dans le cas général
 - **Mais** algorithme de Baum-Welch : «gradient-like»
 - Inévitablement piégé dans des optima locaux \implies métaheuristiques
- Difficultés du problème :
 - Optimisation continue
 - Avec contraintes de stochasticité
 - Un polynôme de degré élevé
 - Beaucoup de variables
 - Nombreux optima locaux
 - Nombreuses «symétries» dans l'espace de recherche (renumérotation des états)
 - Vitesse de convergence de Baum-Welch imprévisible

Les trois problèmes fondamentaux des MMC III

- Les heuristiques et méta heuristiques pour l'apprentissage de MMC avec le critère de maximum de vraisemblance :

Algorithme	
Baum-Welch	(Baum et Eagon, 1967)
Recuit simulé	(Paul, 1985), (Haman et Al Ani, 1996)
Recherche tabou	(Chen et al, 2004)
Algorithme génétique	(Slimane et al, 1999), (Brouard, 1999), (Thomsen, 2002)
Apprentissage incrémental à base de population	(Maxwell et Anderson, 1999)
API (fourmis artificielles)	(Monmarché, 2000), (Aupetit, 2005)
Optimisation par essaim particulaire	(Rasmussen, 2003), (Aupetit, 2005)

L'apprentissage de MMC avec l'OEP

Constats

- Problème d'optimisation **continue** :
 - ⇒ l'OEP est une approche pertinente
- L'algorithme de Baum-Welch converge vers des **optima locaux** à des **vitesses imprévisibles**
 - ⇒ hybrider l'OEP et n'utiliser que quelques itérations de Baum-Welch sur chaque solution
- Optimisation sous contraintes :
 - 1 Utiliser une fonction de pénalité : possible mais pas idéal
 - 2 Ignorer les contraintes puis «normaliser» lorsque nécessaire :
 - Principe suivi par (Rasmussen, 2003)
 - **Mais** ajoute une infinité d'optima locaux et de «symétries»
 - 3 Incorporer naturellement les contraintes en changeant d'espace de solutions

Constats

- Problème d'optimisation **continue** :
 - ⇒ l'OEP est une approche pertinente
- L'algorithme de Baum-Welch converge vers des **optima locaux** à des **vitesses imprévisibles**
 - ⇒ hybrider l'OEP et n'utiliser que quelques itérations de Baum-Welch sur chaque solution
- Optimisation sous contraintes :
 - 1 Utiliser une fonction de pénalité : possible mais pas idéal
 - 2 Ignorer les contraintes puis «normaliser» lorsque nécessaire :
 - Principe suivi par (Rasmussen, 2003)
 - **Mais** ajoute une infinité d'optima locaux et de «symétries»
 - 3 **Incorporer naturellement les contraintes en changeant d'espace de solutions**

L'espace vectoriel des MMC

- Trois étapes :

- ① L'espace des vecteurs stochastiques comme espace vectoriel :

Soit $P_K^* = \{(x_1, \dots, x_K) / \forall i = 1..K, x_i > 0 \text{ et } \sum_{i=1}^K x_i = 1\}$

Soient $\oplus : P_K^* \rightarrow P_K^*$ et $\odot : \mathbb{R} \times P_K^* \rightarrow P_K^*$ tels que

$$(x \oplus y)_i = \frac{x_i y_i}{\sum_{k=1}^K x_k y_k} \text{ et } (c \odot y)_i = \frac{x_i^c}{\sum_{k=1}^K x_k^c}$$

(P_K^*, \oplus, \odot) est un espace vectoriel

- ② L'espace des matrices stochastiques comme espace vectoriel :

$$M_{X \times Y}^* \equiv P_Y^* \times \dots \times P_X^*$$

- ③ L'espace des MMC comme espace vectoriel :

$$\Lambda^* = M_{1 \times N}^* \times M_{N \times N}^* \times M_{N \times M}^*$$

$(\Lambda^*, \oplus, \odot)$ espace vectoriel

- Tout vecteur de Λ^* est un MMC et vice versa (sauf ...)

L'OEP simple sur Λ^*

- La structure vectorielle $(\Lambda^*, \oplus, \odot)$ permet une adaptation facile de l'OEP
- Les équations de mouvement vectorielles :

- 1 Déplacement :

$$\mathbf{x}_i(t+1) = BW(\mathbf{x}_i(t) \oplus \mathbf{v}_i(t))$$

- 2 Compensation de l'opérateur d'optimisation local :

$$\mathbf{v}_i(t) = \mathbf{x}_i(t+1) \ominus \mathbf{x}_i(t)$$

- 3 Mise à jour de la vitesse :

$$\begin{aligned}\mathbf{v}_i(t+1) &= \boldsymbol{\omega} \odot \mathbf{v}_i(t) \\ &\quad \oplus (\mathbf{c}_1 \cdot \mathcal{U}_1) \odot (\mathbf{x}_i^+(t) \ominus \mathbf{x}_i(t)) \\ &\quad \oplus (\mathbf{c}_2 \cdot \mathcal{U}_2) \odot (\hat{\mathbf{x}}_i(t) \ominus \mathbf{x}_i(t))\end{aligned}$$

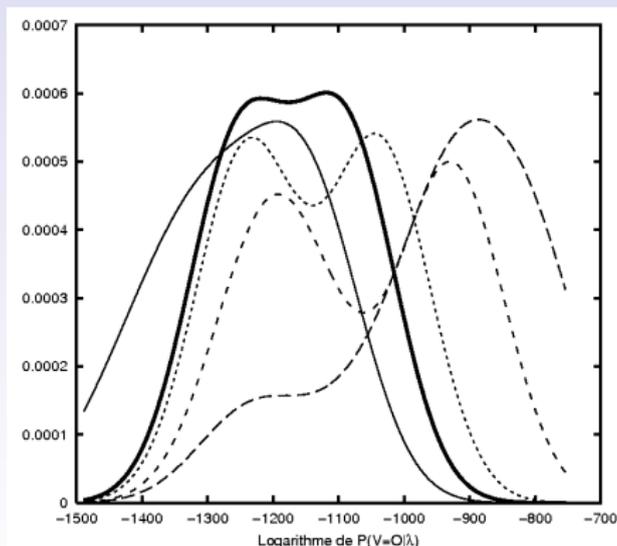
Etudes expérimentales

Recherche de bons paramètres I

- Apprentissage d'un jeu de 8 images :
 - Plus de 48 000 combinaisons de paramètres
 - 15 apprentissages par combinaison
 - 1000 itérations de Baum-Welch ou 30000 évaluations de MMC par apprentissage
 - 10, 20 et 40 états cachés
 - Voisinage circulaire de taille V
- Bon paramètre :
 - Garantit de bonnes performances quasi-indépendamment des autres paramètres
 - Pas de garantie d'optimalité mais relativement indépendant de l'observation apprise

Recherche de bons paramètres II

- Recherche manuelle des bons paramètres grâce à des courbes moyennes estimant de façon lissée la probabilité d'apparition d'un paramètre en fonction de la performance :



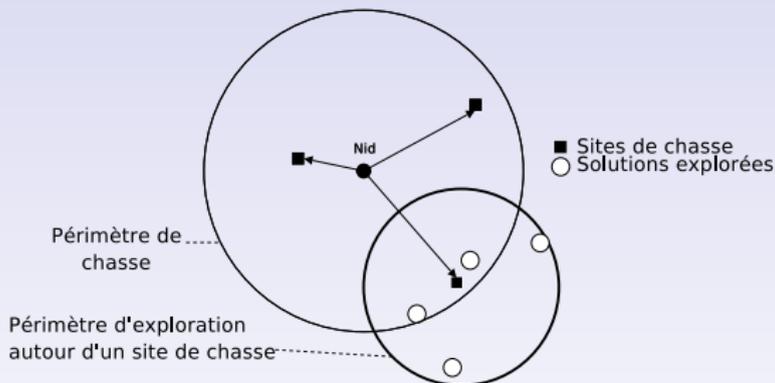
Recherche de bons paramètres III

- Résultats :

	\mathcal{N}_{BW}	\mathcal{N}	ω	c_1	c_2	V
OEP1	0	20	0.8	1.5	2.5	9
OEP2	2	30	0.8	1.5	0	0
OEP3	5	50	0.4	0	0.5	6

Comparaison des performances

- L'algorithme de colonie de fourmis API :
 - Inspiré du comportement des *Pachycondyla Apicalis*



- Adapté :
 - sur la structure vectorielle (Aupetit, 2005),
 - sans structure vectorielle (Monmarché, 2000)
- Les mêmes principes de réglage ont été appliqués sur API

Résultats I

$\bar{e}(A)$ performance moyenne normalisée sur [0;1] sur plus d'images et pour 5000 évaluations de MMC

Algorithme A	Moyenne $\bar{e}(A)$	\mathcal{N}_{BW}
Random0	14.27%	0
OEP1	59.23%	0
Random2	67.94%	2
Random5	82.32%	5
API vectoriel	98.89%	2
OEP3	99.22%	5
API ordinaire	99.46%	5
OEP2	99.53%	2

- Sans hybridation : performances basses, coût de calcul plus élevé
- Avec hybridation : performances élevées
- L'OEP est compétitive avec le plus performant de nos algorithmes

Résultats II

- Evolution du classement selon le nombre total d'itérations de Baum-Welch

Algorithme A	100 itérations		1000 itérations		5000 itérations	
	$\bar{e}(A)$	Ecart	$\bar{e}(A)$	Ecart	$\bar{e}(A)$	Ecart
Random2	53.57%	53%	0.00%	0%	67.94%	8%
Random5	81.67%	82%	50.26%	8%	82.32%	7%
API vectoriel	3.48%	10%	92.71%	3%	98.89%	1%
OEP3	N/A	N/A	84.18%	6%	99.22%	1%
API ordinaire	98.82%	1%	98.24%	2%	99.46%	1%
OEP2	64.00%	16%	99.34%	1%	99.53%	1%

Conclusion

Conclusion

- L'adaptation de l'OEP au problème de l'apprentissage de MMC se fait simplement :
 - $(\Lambda^*, \oplus, \odot)$
 - Hybridation avec l'algorithme de Baum-Welch
- L'hybridation est rentable
- Compétitif avec nos algorithmes les plus efficaces
- Efficacité inférieure avec peu d'itérations de Baum-Welch mais se rattrape par la suite

Perspectives

- L'OEP que nous avons considérée est très simple
- Les équations de mouvement vectorielles ne sont pas les plus efficaces pour l'OEP :
 - D'autres équations de mouvement permettraient peut être d'améliorer les performances ?
- La structure d'espace $(\Lambda^*, \oplus, \odot)$ pourrait être utilisée :
 - pour mieux comprendre ce que font nos méta heuristiques d'apprentissage ?
 - pour adapter naturellement d'autres méta heuristiques à ce problème ?
- Nos algorithmes se concentrent sur le maximum de vraisemblance :
 - Sont-ils performants pour les autres critères ?