

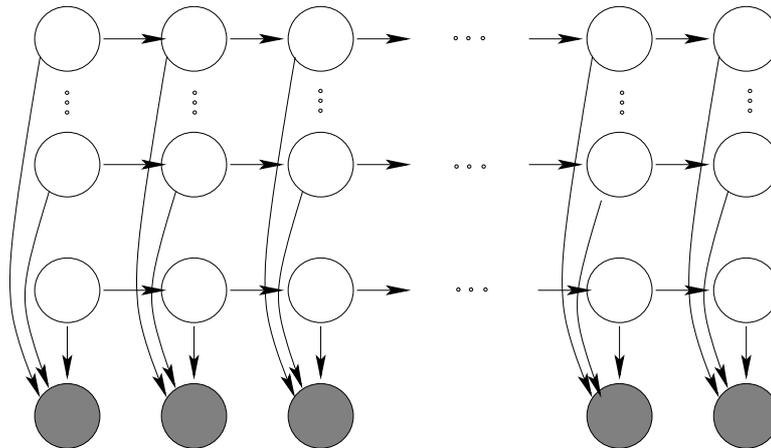
# Mastere M2 MVA 2008 - Modèles graphiques

Exercices à rendre pour le 17 décembre 2008.

Ces exercices doivent s'effectuer seul.

## 1 Modèles de Markov cachés factoriels

On considère le modèle de Markov caché factoriel suivant, avec  $M$  chaînes de Markov cachées de longueur  $T$  (les variables grisées sont toujours observées) :



1. Quelle est la complexité de l'algorithme de l'arbre de jonction pour l'inférence? (on supposera que la  $i$ -ième chaîne a  $r_i$  états).
2. Quelle serait la complexité en utilisant (naivement) une seule chaîne de Markov avec  $\prod_i r_i$  états? (si  $x_1^t, \dots, x_M^t$  sont les états des  $M$  chaînes à l'instant  $t$ , l'état de cette chaîne est  $x^t = (x_1^t, \dots, x_M^t)$ ).

## 2 Modèles orientés Gaussiens

Soient  $n$  variables aléatoires  $X_1, \dots, X_n$  de densité jointe Gaussienne de moyenne  $\mu \in \mathbb{R}^n$  et de matrice de covariance  $\Sigma \in \mathbb{R}^{n \times n}$ .

1. Rappeler (sans démonstration) la condition nécessaire et suffisante pour que la loi Gaussienne  $\mathcal{N}(\mu, \Sigma)$  se factorise dans un graphe non orienté. Le but de cet exercice est de déterminer une caractérisation similaire pour les graphes orientés.
2. On suppose que la loi  $\mathcal{N}(\mu, \Sigma)$  se factorise dans un graphe orienté, caractérisé par les parents  $\pi_i$  de chaque sommet  $i$ ,  $i = 1, \dots, n$ . Montrer que la loi de  $X_i$

sachant  $X_{\pi_i} = x_{\pi_i}$  est une Gaussienne et exprimer sa moyenne et sa variance en fonction de  $x_{\pi_i}$ ,  $\mu$  et  $\Sigma$ .

3. Etant donnée une matrice symétrique définie positive  $M$ , la décomposition de Cholesky est l'unique factorisation  $M = RR^\top$  où  $R$  est triangulaire inférieure avec des éléments positifs sur la diagonale. La décomposition  $LDL^\top$  est l'unique factorisation  $M = LDL^\top$  où  $D$  est diagonale avec diagonale strictement positive et  $L$  triangulaire inférieure avec diagonale constante égale à 1 (on passe de l'une à l'autre par  $R = LD^{1/2}$ , où  $D^{1/2}$  est la matrice diagonal des racines carrées positives des éléments diagonaux de  $D$ ). Dans cet exercice, on considère les décompositions de  $M = \Sigma^{-1}$ .

On suppose que l'ordre  $(n, \dots, 1)$  est topologique pour un graphe orienté donné  $G = (V, E)$  (défini par l'ensemble des parents de chaque noeud), i.e.,  $\pi_i \subset \{i + 1, \dots, n\}$ . Montrer que la loi  $\mathcal{N}(\mu, \Sigma)$  se factorise dans le graphe orienté ssi

$$\forall i \in \{1, \dots, n\}, (j \neq i \text{ et } j \notin \pi_i) \Rightarrow R_{ji} = 0.$$

ou, de manière équivalente,

$$\forall i \in \{1, \dots, n\}, (j \neq i \text{ et } j \notin \pi_i) \Rightarrow L_{ji} = 0.$$

(i.e., les modèles graphiques orientés correspondent à une factorisation parcimonieuse de  $\Sigma^{-1}$ ).

INDICATIONS : (a) se ramèner à  $\mu = 0$ , (b) exprimer  $p(x) (= \prod_{i=1}^n p(x_i | x_{i+1}, \dots, x_n))$  à l'aide de  $L$  et  $D$ , (c) utiliser la caractérisation suivante : si  $(n, \dots, 1)$  est un ordre topologique pour  $G$ , alors  $p(x)$  se factorise dans  $G$  ssi  $\forall i < n, p(x_i | x_{i+1}, \dots, x_n) = p(x_i | x_{\pi_i})$ .

4. Montrer comment les résultats de la question 2.2 permettent de calculer directement  $L$  à partir de  $\Sigma$ , sous les hypothèses que la loi  $\mathcal{N}(\mu, \Sigma)$  se factorise et que l'ordre  $(n, \dots, 1)$  est topologique.
5. BONUS : si la factorisation  $LDL^\top$  de  $\Sigma^{-1}$  est connue, quelle est la complexité numérique des produits "matrice-vecteur" de la forme  $\Sigma y$  et  $\Sigma^{-1}y$ , pour  $y \in \mathbb{R}^n$  donné.
6. "DOUBLE" BONUS : si la loi  $\mathcal{N}(\mu, \Sigma)$  se factorise dans un modèle graphique non orienté, et que  $\Sigma^{-1}$  est donnée, quelle est la complexité numérique des produits "matrice-vecteur" de la forme  $\Sigma y$  et  $\Sigma^{-1}y$ , pour  $y \in \mathbb{R}^n$  donné (on supposera connu la largeur arborescente  $t$  du graphe ainsi qu'un ordre d'élimination associé).

### 3 Apprentissage de la structure d'un arbre

Dans cet exercice, un algorithme pour apprendre la structure d'un modèle graphique à partir de données i.i.d sera dérivé pour les arbres.

1. Préliminaire I (entropie marginale) : Soit une variable aléatoire discrète  $X$  à valeurs dans un ensemble fini  $\mathcal{X}$ . Soit  $\eta(x) = p(X = x)$  le vecteur de paramètres. Soit un échantillon i.i.d  $(x^n)$ ,  $n = 1, \dots, N$  de taille  $N$  de cette variable. On note  $\hat{p}(x)$  la densité empirique, définie par  $\hat{p}(x) = \frac{1}{N} \sum_{n=1}^N \delta(x^n = x)$ . Montrer que la log-vraisemblance conditionnelle maximale  $\max_{\eta} \sum_{n=1}^N \log p(x^n | \eta)$  est égale à :

$$\max_{\eta} \sum_{n=1}^N \log p(x^n | \eta) = -NH(X),$$

où  $H(X)$  est l'entropie empirique de  $X$ , définie par  $H(X) = - \sum_{x \in \mathcal{X}} \hat{p}(x) \log \hat{p}(x)$ .

2. Préliminaire II (entropies jointe et conditionnelle) : Soient deux variables aléatoires discrètes  $X, Y$  à valeurs dans des ensembles finis  $\mathcal{X}$  et  $\mathcal{Y}$ . Soit  $\eta(x, y) = p(Y = x | X = x)$  la matrice de paramètres de la loi conditionnelle. Soit un échantillon i.i.d  $(x^n, y^n)$ ,  $n = 1, \dots, N$  de taille  $N$  de ces deux variables. On note  $\hat{p}(x, y)$  la densité empirique, définie par  $\hat{p}(x, y) = \frac{1}{N} \sum_{n=1}^N \delta(x^n = x) \delta(y^n = y)$ . Montrer que la log-vraisemblance conditionnelle maximale  $\max_{\eta} \sum_{n=1}^N \log p(y^n | x^n, \eta)$  est égale à :

$$\max_{\eta} \sum_{n=1}^N \log p(y^n | x^n, \eta) = N(H(X) - H(X, Y))$$

où  $H(X, Y)$  l'entropie (jointe) empirique de  $(X, Y)$  est définie par

$$H(X, Y) = - \sum_{y \in \mathcal{Y}} \sum_{x \in \mathcal{X}} \hat{p}(x, y) \log \hat{p}(x, y).$$

3. On considère maintenant  $P$  variables aléatoires  $X_1, \dots, X_P$  à supports finis  $\mathcal{X}_1, \dots, \mathcal{X}_P$ . On considère  $N$  observations i.i.d. de ces  $P$  variables,  $(x_p^n)$ ,  $p = 1, \dots, P$ ,  $n = 1, \dots, N$ . On note  $\hat{p}(x_1, \dots, x_P)$  la densité empirique, définie comme suit :

$$\hat{p}(x_1, \dots, x_P) = \frac{1}{N} \sum_{n=1}^N \delta(x_1^n = x_1) \cdots \delta(x_P^n = x_P).$$

Cette densité jointe permet de définir des densités marginales  $\hat{p}(x_p, x_q)$  et  $\hat{p}(x_q)$  par marginalisation.

Soit un arbre couvrant orienté  $T$  à  $P$  sommets (i.e., un DAG connexe avec au plus un parent par sommet). Quelle est la paramétrisation la plus générale pour une loi  $p(x_1, \dots, x_P)$  se factorisant dans un tel DAG ? Montrer qu'une fois maximisée par rapport à ces paramètres, la log-vraisemblance des données est égale à :

$$\ell(T) = N \sum_{p=1}^P \{H(X_{\pi_p(T)}) - H(X_p, X_{\pi_p(T)})\}$$

(avec les conventions  $H(X_p, X_{\emptyset}) = H(X_p)$  et  $H(X_{\emptyset}) = 0$ )

4. Pour tout  $p, q$ , on appelle information mutuelle empirique la quantité  $I(X_p, X_q) = -H(X_p, X_q) + H(X_p) + H(X_q)$ . Exprimer cette quantité comme une divergence de Kullback-Leibler et montrer qu'elle est positive ou nulle.
5. Exprimer  $\ell(T)$  à l'aide des informations mutuelles. Comment maximiser  $\ell(T)$  par rapport à l'arbre  $T$ ? Décrire un algorithme permettant d'apprendre la structure de l'arbre ayant le maximum de vraisemblance. Quelle est sa complexité?

## 4 Implémentation - HMM

On considère les mêmes données d'apprentissage que le devoir précédent, dans le fichier "EMGaussienne.dat", mais cette fois-ci en considérant la structure temporelle, i.e., les données sont de la forme  $u_t = (x_t, y_t)$  où  $u_t = (x_t, y_t) \in \mathbb{R}^2$ , pour  $t = 1, \dots, T$ . Le but de cet exercice est d'implémenter l'inférence dans les HMM ainsi que l'algorithme EM pour l'apprentissage des paramètres. Il est conseillé d'utiliser le code du devoir précédent.

On considère le modèle HMM suivant avec une chaîne  $(q_t)$  à  $K=4$  états et matrice de transition  $a \in \mathbb{R}^{4 \times 4}$ , et des "probabilités d'émission" Gaussiennes :  $u_t | q_t = i \sim \mathcal{N}(\mu_i, \Sigma_i)$ .

1. Implémenter les récursions  $\alpha$  et  $\beta$  vues en cours et dans le polycopié pour estimer  $p(q_t | u_1, \dots, u_T)$  et  $p(q_t, q_{t+1} | u_1, \dots, u_T)$ .
2. Calculer les équations d'estimation de EM.
3. Implémenter l'algorithme EM pour l'apprentissage (on pourra initialiser les moyennes et les covariances avec celles trouvées dans le devoir précédent).
4. Implémenter l'inférence pour estimer la séquence d'états la plus probable, i.e.  $\arg \max_q p(q_1, \dots, q_T | y_1, \dots, y_T)$ , et représenter le résultat obtenu avec les données (pour le jeu de paramètres appris par EM).
5. Commenter les différents résultats obtenus avec ceux du devoir précédent. En particulier, comparer les log-vraisemblances, sur les données d'apprentissage, ainsi que sur les données de test (dans "EMGaussienne.test").