

Sur l'apprentissage de Réseaux Bayésiens à partir de bases d'exemples incomplètes et application à la classification

Olivier FRANÇOIS et Philippe LERAY,
Laboratoire **LITIS**, Rouen.

Rencontres Inter-Associations

La classification et ses applications

AFIA - ARIA - EGC - INFORSID - SFC - SFDS - LMO - ASTI

21 Mars 2006

Plan

- 1 *Introduction*
 - Données manquantes et entrées incomplètes
 - Classification : le Réseau Bayésien Naïf
- 2 *Espérance-Maximisation Structurale*
 - La méthode AMS-EM
 - Notre adaptation : MWST-EM
- 3 *Le réseau Bayésien Naïf Augmenté par un Arbre*
 - Apprendre une Structure pour la Classification à partir de Données Incomplètes
 - Résultats
- 4 *Conclusions et Perspectives*
References 15

Problématique

Soit \mathcal{X} un système complexe.

\mathcal{X} est représenté par de nombreux attributs $\{X_i\}_{1 \leq i \leq n}$.

- Certains attributs sont observés **systematiquement**,
- d'autres sont observés **occasionnellement**,
 - état *critique* du système ?
 - mesure *couteuse* ?...
- et de nombreux autres ne sont **jamais** observés,
 - parce que leur *influence/pertinence est faible* ?
 - parce que l'on ne *pas connaissance* de leur intérêt ?...

Par Exemple : Pour une base de 2000 exemples sur 20 attributs,

20% des mesures sont manquantes complètement au hasard

⇒ en moyenne *seulement* 23 cas complets (c-à-d %EI \simeq 99%)

Types de données manquantes

Notations :

$$\mathbf{D} = \langle \mathbf{O}, \mathbf{H} \rangle = ((d_{ij}))_{n \times m}$$

$\mathbf{R} = ((r_{ij}))_{n \times m}$, une matrice où $r_{ij} = 1$ si d_{ij} est manquant, 0 sinon.

θ , paramètres de la loi qui a généré \mathbf{D} ,

μ , paramètres de la loi qui a généré \mathbf{R} .

Données manquantes ?

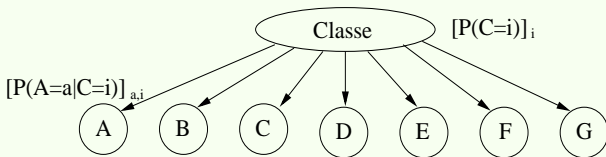
Rubin (1976)

$$\mathbb{P}(\mathbf{O}, \mathbf{H}, \mathbf{R} | \theta, \mu) = \mathbb{P}(\mathbf{O}, \mathbf{H} | \theta) \times \mathbb{P}(\mathbf{R} | \mathbf{O}, \mathbf{H}, \mu)$$

- MCAR : $\mathbb{P}(\mathbf{R} | \mathbf{O}, \mathbf{H}, \mu) = \mathbb{P}(\mathbf{R} | \mu)$
- MAR : $\mathbb{P}(\mathbf{R} | \mathbf{O}, \mathbf{H}, \mu) = \mathbb{P}(\mathbf{R} | \mathbf{O}, \mu)$
- NMAR : $\mathbb{P}(\mathbf{R} | \mathbf{O}, \mathbf{H}, \mu)$, cas non ignorables.

Le réseau Naïf

Supposons que *la classe a une influence sur toutes les variables, mais* indépendamment

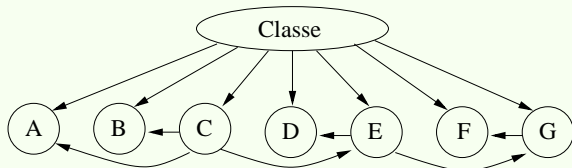


Ce qu'il est possible de faire :

- l'apprentissage des paramètres avec DI (par ex. avec EM),
- l'inférence avec des Données Incomplètes.

Le réseau Naïf

Supposons que *la classe a une influence sur toutes les variables, mais plus* indépendamment



Ce qu'il est *toujours* possible de faire :

- l'apprentissage des paramètres avec DI (par ex. avec EM),
- l'inférence avec des Données Incomplètes.

et

si l'on veut ajouter des dépendances automatiquement : ?

(peu de méthodes efficaces à partir de D

Donner un score à partir d'une base incomplète

Soit $\mathcal{S}(\mathcal{M}|\mathbf{D}_C)$, score pour un modèle \mathcal{M} et
des données complètes \mathbf{D}_C .

→ approximation de \mathcal{S} pour \mathcal{M} et la base incomplète $\mathbf{D} = \langle \mathbf{O}, \mathbf{H} \rangle$

$$\mathcal{Q}^{\mathcal{S}}(\mathcal{M}|\mathbf{D}) = \mathbb{E}_{\mathbf{H} \sim \mathbb{P}(\mathbf{H}|\mathbf{O}, \mu)} [\mathcal{S}(\mathcal{M}|\mathbf{O}, \mathbf{H})]$$

Mais la loi $\mathbb{P}(\mathbf{H}|\mathbf{O}, \mu)$ est inconnue.

Principe EM

Supposons que le modèle \mathcal{M}^0 a généré la base \mathbf{D} alors

$$\begin{aligned} \mathcal{Q}^{\mathcal{S}}(\mathcal{M}|\mathbf{D}) &\approx \mathcal{Q}^{\mathcal{S}}(\mathcal{M} : \mathcal{M}^0|\mathbf{D}) = \mathbb{E}_{\mathbf{H} \sim \mathbb{P}(\mathbf{H}|\mathbf{O}, \mathcal{M}^0)} [\mathcal{S}(\mathcal{M}|\mathbf{O}, \mathbf{H})] \\ \mathcal{Q}^{\mathcal{S}}(\mathcal{M} : \mathcal{M}^0|\mathbf{D}) &= \sum_{\mathbf{H}} \mathcal{S}(\mathcal{M}|\mathbf{O}, \mathbf{H}) \mathbb{P}(\mathbf{H}|\mathbf{O}, \mathcal{M}^0) \end{aligned}$$

Ou la loi *a posteriori* $\mathbb{P}(\mathbf{H}|\mathbf{O}, \mathcal{M}^0)$ est connue.

Structural-EM

- Choisir un modèle \mathcal{M}^0 ($\Rightarrow \mathbb{P}(\mathbf{H}|\mathcal{M}^0)$)
- Trouver un modèle \mathcal{M}^{i+1} qui maximise** un score $Q^S(\mathcal{M} : \mathcal{M}^i | \mathbf{D})$
- Utiliser le nouveau modèle comme base pour l'itération suivante jusqu'à convergence.

**EM généralisé : augmente le score

AMS-EM : le nouveau modèle est choisi parmi les *voisins* du graphe courant Friedman (1997)

MWST-EM

AMS-EM : le nouveau modèle est choisi parmi les *voisins* du graphe courant Friedman (1997) → nombreuses itérations

MWST-EM : nous trouvons le 'meilleur' modèle dans l'espace des arbres Leray & François (2005) → peu d'itérations

Pour cela, on utilise un algorithme de type Kruskal sur la matrice de score suivante :

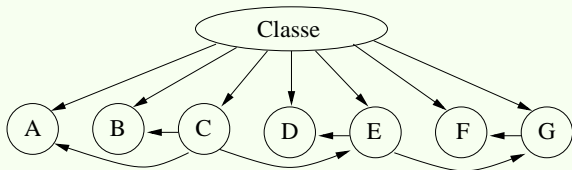
$$\left[M_{ij}^Q \right]_{i,j} = \left[Q^{bic}(X_i, P_i = \{X_j\}, \Theta_{X_i|X_j} : T^*, \Theta^*) - Q^{bic}(X_i, P_i = \emptyset, \Theta_{X_i} : T^*, \Theta^*) \right]$$

$$\text{où } Q^{BIC}(\mathcal{G}, \theta : \mathcal{G}^*, \Theta^*) = \sum_i Q^{bic}(X_i, P_i, \Theta_{X_i|P_i} : \mathcal{G}^*, \Theta^*) \quad \text{et}$$

$$Q^{bic}(X_i, P_i, \Theta_{X_i|P_i} : \mathcal{G}^*, \Theta^*) = \sum_{X_i=X_k} \sum_{P_i=pq_j} N_{ijk}^* \log \theta_{ijk} - \frac{\log N}{2} \text{Dim}(\Theta_{X_i|P_i})$$

$$\text{avec } N_{ijk}^* = E_{\mathcal{G}^*, \Theta^*} [N_{ijk}] = N * P(X_i = x_k, P_i = pq_j | \mathcal{G}^*, \Theta^*).$$

TAN-EM



Pour augmenter le réseau naïf par un arbre,

la classe fait toujours partie de l'ensemble des parents

et la matrice de score devient ($i, j \neq \text{classe}$) :

$$\left[M_{ij}^Q \right]_{i,j} = \left[Q^{bic}(X_i, P_i = \{C, X_j\}, \theta_{X_i|X_j} : \mathcal{T}^*, \Theta^*) \right. \\ \left. - Q^{bic}(X_i, P_i = \{C\}, \theta_{X_i} : \mathcal{T}^*, \Theta^*) \right]$$

Résultats

	N	N app	N test	#C	%EI
Hepatitis	20	90	65	2	8.4
House	17	290	145	2	46.7
Horse	28	300	300	2	88.0
Thyroid	22	2800	972	2	29.9
Mushrooms	23	5416	2708	2	30.5

	NB-EM	MWST-EM	TAN-EM	AMS-EM	AMS-EM+T
Hepatitis	70.8% -1224 ; 29	73.8% -1147 ; 90	75.4% -1148 ; 88	66.1% -1211.5 ; 1213	66.1% -1207 ; 1478
House	89.7% -2203 ; 110	93.8% -2518 ; 157	92.4% -2022 ; 180	92.4% -2524 ; 1732	93.8% -2195 ; 3327
Horse	75% -5589 ; 227	77.9% -5199 ; 656	80.9% -5354 ; 582	66.2% -5348 ; 31807	66.2% -5318 ; 10054
Thyroid	95.3% -39348 ; 1305	93.8% -38881 ; 3173	96.2% -38350 ; 3471	93.8% -38303 ; 17197	93.8% -39749
Mushrooms	92.8% -97854 ; 2028	74.7% -108011 ; 6228	91.3% -87556 ; 5987	74.9% -111484 ; 70494	74.9% -110828 ; 50753

Conclusions

La méthode TAN-EM permet d'obtenir :

- de bonnes performances en classification,
- de bonnes vraisemblances des modèles obtenus,
- un excellent rapport *performances/rapidité* (car basée sur le réseau naïf et MWST-EM Leray & François (2005)).

Néanmoins, cette méthode

- est limitée au tâches de classification,
- augmente le NB *forcement* par un arbre.

Perspectives

Pour TAN-EM :

- adaptation à la classification non-supervisée,
- tests sur des données générées (MAR),
- tests en non-supervisé.

Pour MWST-EM :

- passer de l'espace des arbres à
 - l'espace des *forets* (\implies FAN-EM),
 - l'espace des *équivalents de Markov* (\implies GES-EM et BNAN-EM ?),
- remplacer les principes de EM par d'autres (Robust Bayesian Estimator...) \implies NMAR.

Merci pour votre attention.

- Questions ?
- Remarques ?
- Suggestions ?

Friedman, N. (1997).

Learning belief networks in the presence of missing values and hidden variables.

In in the Proceedings of the 14th International Conference on Machine Learning, (pp. 125–133). Morgan Kaufmann.

Leray, P. & François, O. (2005).

Bayesian network structural learning and incomplete data.

In in the proceedings of the International and Interdisciplinary Conference on Adaptive Knowledge Representation and Reasoning (AKRR 2005), Espoo, Finland.

Rubin, D. (1976).

Inference and missing data.

Biometrika, 63, 581–592.