

# Apprentissage de structure à partir de données incomplètes et application à la classification

Olivier François, Philippe Leray  
Francois.Olivier.C.H@gmail.com,  
Philippe.Leray@insa-rouen.fr

Laboratoire d'Informatique, Traitement de l'Information, Systèmes.

3ème Journées Francophones sur les Réseaux bayésiens  
organisées à valenciennes (Nord)

les 18 et 19 juin.



- 1 Introduction**
  - Données manquantes et entrées incomplètes
  - Classification : le Réseau Bayésien Naïf
- 2 Espérance-Maximisation Structurale**
  - La méthode AMS-EM
  - Notre adaptation : MWST-EM
  - Résultats
- 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**

# 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.

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

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

Données manquantes ?

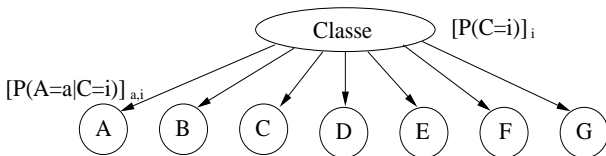
(Rubin, 76)

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

- 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

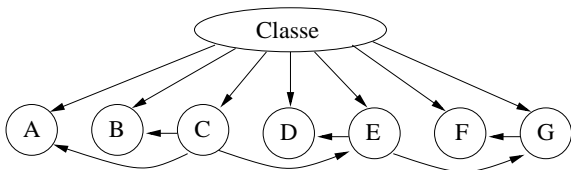


Ce qu'il est possible de faire :

- 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 DI)

# Plan

- 1 Introduction
- 2 Espérance-Maximisation Structurale**
- 3 Le réseau Bayésien Naïf Augmenté par un Arbre
- 4 Conclusions et Perspectives

## Calculer un score

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  $\mathbf{D} = \langle \mathbf{O}, \mathbf{H} \rangle$

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

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

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

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



# Structural-EM

- Choisir un modèle  $\mathcal{M}^0$  ( $\Rightarrow \mathbb{P}(\mathbf{H}|\mathbf{O}, \mathcal{M}^0)$ )
- 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,97)

AMS-EM : le nouveau modèle est choisi parmi les *voisins* du graphe courant (Friedman,97) → nombreuses itérations

MWST-EM : nous trouvons le 'meilleur' modèle dans l'espace des arbres (François & Leray,05) → peu d'itérations

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

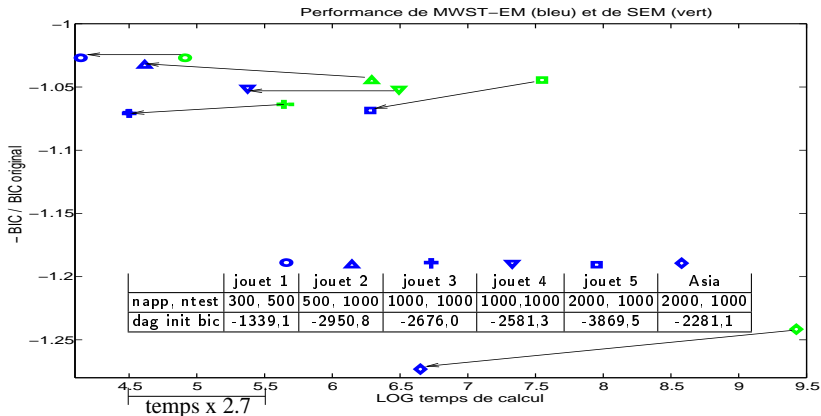
$$\left[ M_{ij}^Q \right]_{i,j} = \left[ Q^{bic}(X_i, P_i = \{X_j\}, \Theta_{X_i|X_j} : T^*, \Theta^*) - Q^{bic}(X_i, P_i = \emptyset, \Theta_{X_i} : T^*, \Theta^*) \right]$$

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

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

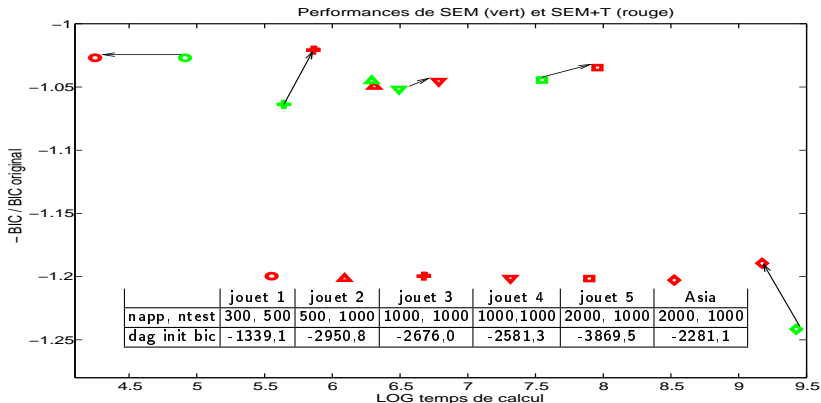
avec  $N_{ijk}^* = E_{\mathcal{G}^*, \Theta^*} [N_{ijk}] = N * P(X_i = x_k, P_i = p_{a_j} | \mathcal{G}^*, \Theta^*)$ .

## Résultats Préliminaires : données MCAR



- MWST-EM donne de 'bons' résultats.
  - L'arbre est proche de la solution
  - Le temps de calcul est faible

# Résultats Préliminaires : données MCAR

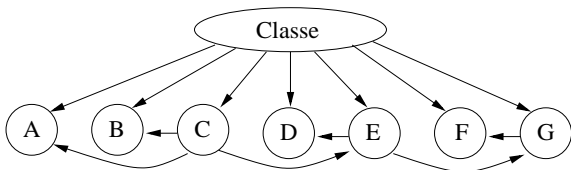


- SEM n'est pas *stable* (initialisation).
- Initialiser SEM avec un arbre optimal
  - stabilise la méthode,
  - *meilleurs résultats* et/ou mêmes résultats plus *rapidement*.

# Plan

- 1 Introduction
- 2 Espérance-Maximisation Structurale
- 3 Le réseau Bayésien Naïf Augmenté par un Arbre**
- 4 Conclusions et Perspectives

## 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 ( $i, j \neq classe$ ) :

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

## 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% <b>-1147</b> ; 90	<b>75.4%</b> <b>-1148</b> ; 88	66.1% -1211.5 ; 1213	66.1% -1207 ; 1478
House	89.7% -2203 ; 110	<b>93.8%</b> -2518 ; 157	92.4% <b>-2022</b> ; 180	92.4% -2524 ; 1732	<b>93.8%</b> -2195 ; 3327
Horse	75% -5589 ; 227	77.9% <b>-5199</b> ; 656	<b>80.9%</b> -5354 ; 582	66.2% -5348 ; 31807	66.2% -5318 ; 10054
Thyroid	95.3% -39348 ; 1305	93.8% -38881 ; 3173	<b>96.2%</b> <b>-38350</b> ; 3471	93.8% <b>-38303</b> ; 17197	93.8% -39749 ; 14482
Mushrooms	<b>92.8%</b> -97854 ; 2028	74.7% -108011 ; 6228	91.3% <b>-87556</b> ; 5987	74.9% -111484 ; 70494	74.9% -110828 ; 59795

# 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 ?).

Néanmoins, cette méthode

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



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 ?

