

cours 3 Apprentissage d'un réseau bayésien

 Master SID — Raisonnement dans l'incertain

©CG(2023)

Généralités sur l'apprentissage

Apprentissage de réseaux bayésiens

- ▶ Objectif : estimer
 - ▶ la structure \mathcal{G} du réseau bayésien
 - ▶ les paramètres $P(X|\text{pa}(X))$ du réseau bayésien
- ▶ En se fondant sur :
 - ▶ une ou plusieurs base(s) de données
 - ▶ complètes ou avec données manquantes
 - ▶ des connaissances *a priori*
 - ▶ contraintes sur la structure du RB
 - ▶ *A priori* sur les paramètres $P(X|\text{pa}(X))$
 - ▶ connaissances expertes, *etc.*



- 1 Apprendre la structure du RB
- 2 Apprendre les paramètres sachant cette structure

cours 3 Apprentissage d'un réseau bayésien

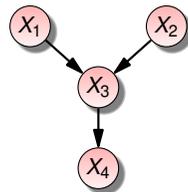
1/27

Apprentissage de paramètres (1/4)

- ▶ Base de données **D complète** : pas de valeur manquante
 $\mathbf{D} = N$ lignes : $\mathbf{D} = \langle d^{(1)}, \dots, d^{(N)} \rangle$
 ligne $d^{(i)}$: instanciation/observation de **toutes** les variables

- ▶ Structure du RB \mathcal{G} connue

X_1	X_2	...	X_n
1	toto	...	0
...
2	titi	...	0



- ▶ Θ : ensemble des paramètres du RB
valeurs des tables $P(X_i|\text{Pa}(X_i))$

Objectif : Estimer Θ qui « colle » le mieux aux données **D**

- ▶ « colle » le mieux \implies le plus vraisemblable

Apprentissage de paramètres (2/4)

Vraisemblance : $\mathcal{L}(\Theta : \mathbf{D}) = P(\mathbf{D}|\Theta)$

- ▶ **Structure \mathcal{G}** \implies indépendances

Estimation par maximum de vraisemblance

- ▶ Θ_i : les paramètres de $P(X_i|\text{Pa}(X_i))$
- ▶ Estimer **indépendamment** chaque Θ_i :
 - ▶ en ne tenant compte que des colonnes X_i et $\text{Pa}(X_i)$ de **D**
 - ▶ en calculant $\text{Argmax}_{\Theta_i} \mathcal{L}(\Theta_i : \mathbf{D})$

cours 3 Apprentissage d'un réseau bayésien

3/27

cours 3 Apprentissage d'un réseau bayésien

2/27

Apprentissage de paramètres (3/4)

Estimation d'une distribution marginale $P(X)$:

- ▶ X : variable aléatoire, domaine : $\Omega_X = \{x_1, \dots, x_n\}$
- ▶ $\Theta = \{\theta_1, \dots, \theta_n\}$ $P(X = x_i | \Theta) = \theta_i$
- ▶ N_i : nombre d'occurrences de x_i dans \mathbf{D}

Théorème

Si $\Theta^* = \text{Argmax}_{\Theta} \mathcal{L}(\Theta : \mathbf{D})$, alors $\theta_i^* = \frac{N_i}{N}$

- ▶ **Démonstration** : Optimum obtenu pour $\frac{\partial \log \mathcal{L}(\Theta : \mathbf{D})}{\partial \theta_i} = 0$
⚠ contrainte : $\sum_{i=1}^n \theta_i = 1$

Exemple d'apprentissage de $P(X)$

- ▶ $X = \ll \text{dé à 6 faces} \gg$ 

- ▶ Base de données \mathbf{D} :

X	1	3	2	4	2	5	1	1	2	3
---	---	---	---	---	---	---	---	---	---	---

- ▶ $N = 10$

- ▶ N_i :

N_1	N_2	N_3	N_4	N_5	N_6
3	3	2	1	1	0

- ▶ Estimation de $P(X = x_j)$:

x_1	x_2	x_3	x_4	x_5	x_6
0,3	0,3	0,2	0,1	0,1	0

Apprentissage de paramètres (4/4)

Estimation de $P(X_i | \text{Pa}(X_i))$ par max de vraisemblance (MLE)

- ▶ r_i : taille du domaine de X_i
domaine de $X_i = \{x_{i1}, \dots, x_{ir_i}\}$
- ▶ q_i : taille du domaine de $\text{Pa}(X_i)$
domaine de $\text{Pa}(X_i) = \{w_{i1}, \dots, w_{iq_i}\}$
- ▶ N_{ijk} : nombre d'occurrences de $(X_i = x_k, \text{Pa}(X_i) = w_{ij})$ dans \mathbf{D}
 $N_{ij} = \sum_k N_{ijk}$
- ▶ $\Theta_i = \{\theta_{ijk} : 1 \leq j \leq q_i, 1 \leq k \leq r_i\}$: paramètres de $P(X_i | \text{Pa}(X_i))$
 $\theta_{ijk} = P(X_i = x_k | \text{Pa}(X_i) = w_{ij}, \Theta_i)$
- ▶ Si $\Theta_i^* = \text{Argmax}_{\Theta_i} \mathcal{L}(\Theta_i : \mathbf{D})$, alors $\theta_{ijk}^* = \frac{N_{ijk}}{N_{ij}}$

⚠ N_{ij} peut être égal à 0 !

Exemple d'apprentissage de $P(X|Y)$

- ▶ $\Omega_X = \{x_1, x_2\}$, $\Omega_Y = \{y_1, y_2\}$

- ▶ Base de données \mathbf{D} :

X	x_1	x_2	x_2	x_2	x_1	x_1	x_1	x_1	x_2	x_2
Y	y_1	y_1	y_1	y_1	y_2	y_2	y_2	y_2	y_2	y_2

- ▶ N_{ijk} N_{ij}

	x_1	x_2	
y_1	1	3	4
y_2	4	2	6

- ▶ Estimation de $P(X|Y)$:

	x_1	x_2
y_1	1/4	3/4
y_2	4/6	2/6

Apprentissage de paramètres avec *a priori* (1/2)

- ▶ **A priori** : distribution $\pi(\Theta)$ sur les paramètres
⇒ Estimation par maximum *a posteriori* (MAP)

$$\Theta^* = \text{Argmax}_{\Theta} P(\Theta|\mathbf{D})$$

- ▶ Formule de Bayes : $P(\Theta|\mathbf{D}) = \frac{P(\mathbf{D}|\Theta) \times \pi(\Theta)}{P(\mathbf{D})}$
 - ▶ $P(\mathbf{D}) = \sum_{\theta} P(\mathbf{D}|\Theta = \theta) \times \pi(\theta) = \text{constante pour le Argmax}$
⇒ $\Theta^* = \text{Argmax}_{\Theta} P(\mathbf{D}|\Theta) \times \pi(\Theta) = \text{Argmax}_{\Theta} \prod_{i=1}^n \mathcal{L}(\Theta_i : \mathbf{D})\pi(\Theta)$
 - ▶ **Hypothèse** : indépendance des paramètres : $\pi(\Theta) = \prod_{i=1}^n \pi(\Theta_i)$
⇒ $\Theta^* = \text{Argmax}_{\Theta} \prod_{i=1}^n \mathcal{L}(\Theta_i : \mathbf{D})\pi(\Theta_i)$
- $$\Rightarrow \Theta_i^* = \text{Argmax}_{\Theta_i} P(\mathbf{D}|\Theta_i) \times \pi(\Theta_i)$$

A priori classique

Distribution de Dirichlet

- ▶ Y : variable de domaine le simplexe k -dimensionnel :
 $\{k\text{-uplets } (y_1, \dots, y_k) \text{ t.q. } y_i \geq 0 \text{ pour tout } i, \text{ et } \sum_{i=1}^k y_i = 1\}$
⇒ $Y = \text{ensemble de distributions de probabilité}$
- ▶ Soit $\alpha = \{\alpha_1, \dots, \alpha_k\}$ t.q. $\alpha_i > 0$ pour tout i
- ▶ $\text{Dir}(Y, \alpha) = \text{distribution de probabilité définie sur } \Omega_Y \text{ par :}$

$$\text{Dir}(Y = y, \alpha) = \frac{1}{B(\alpha)} \prod_{i=1}^k y_i^{\alpha_i - 1}$$

avec $B(\cdot) = \text{constante de normalisation} = \text{fonction Beta}$

- ▶ Justification : Geiger & Heckerman (1997)

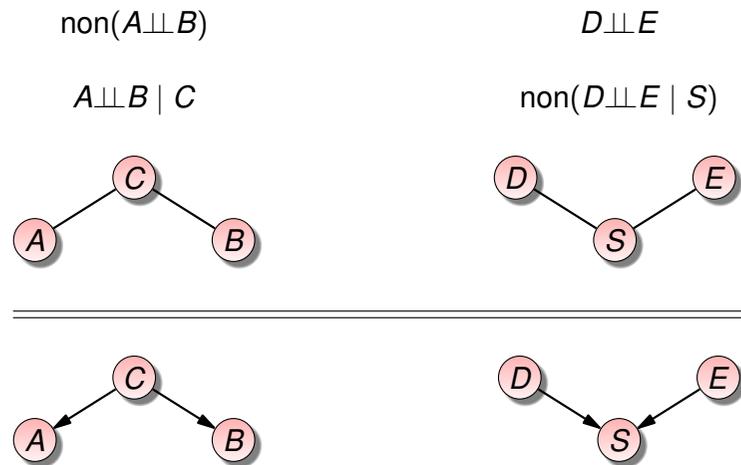
Apprentissage de paramètres avec *a priori* (2/2)

Estimation par Max a posteriori (MAP)

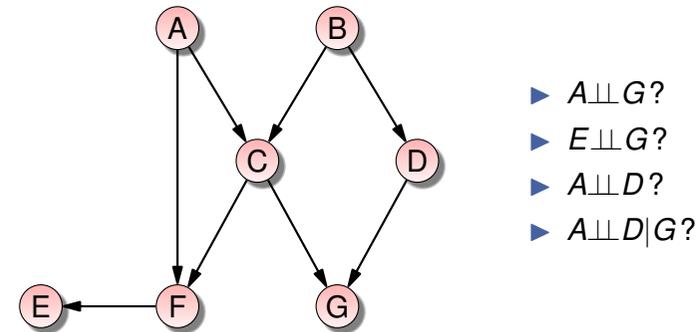
- ▶ *A priori* de Dirichlet d'hyperparamètres α_{ijk}
- ▶ Soit $\alpha_{ij} = \sum_{k=1}^{r_i} \alpha_{ijk}$
- ▶ Si $\Theta_i^* = \text{Argmax}_{\Theta_i} P(\mathbf{D}|\Theta_i) \times \pi(\Theta_i)$ alors $\theta_{ijk}^* = \frac{N_{ijk} + \alpha_{ijk} - 1}{N_{ij} + \alpha_{ij} - r_i}$

② Indépendances et graphe

Rappel de l'épisode précédent



Indépendances et modèle graphique



- ⇒ Raisonnement sur la partie graphique du modèle
- ⇒ Vérifications d'indépendances conditionnelles sans connaître les valeurs des probabilités !

d-séparation

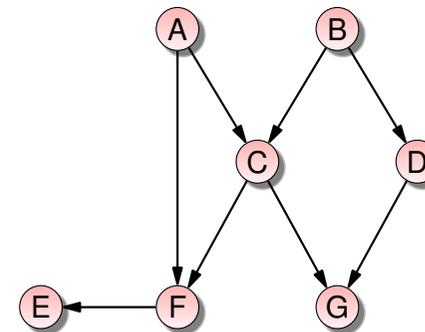
Chaîne $\langle X_1, \dots, X_n \rangle$

- ▶ ensemble de nœuds $\{X_1, \dots, X_n\}$
- ▶ pour tout $i \in \{1, \dots, n-1\}$, le graphe contient l'arc $X_i \rightarrow X_{i+1}$ ou $X_{i+1} \rightarrow X_i$

Chaîne $\langle X_1, \dots, X_n \rangle$ bloquée par un ensemble Z

- ▶ Bloquée si et seulement si $\exists i \in \{2, \dots, n-1\}$ tel que l'une des 2 propriétés ci-dessous est vérifiée :
 - 1 (X_{i-1}, X_i, X_{i+1}) est une V-structure : $X_{i-1} \rightarrow X_i \leftarrow X_{i+1}$, et ni X_i ni ses descendants ne sont dans Z
 - 2 (X_{i-1}, X_i, X_{i+1}) n'est pas une V-structure et $X_i \in Z$

d-séparation – bis



- ▶ Chaîne $\langle D, B, C, A \rangle$ bloquée par \emptyset ?
- ▶ Chaîne $\langle D, B, C, A \rangle$ bloquée par $\{E\}$?
- ▶ Chaîne $\langle D, B, C, F, A \rangle$ bloquée par $\{E\}$?
- ▶ Chaîne $\langle D, G, C, A \rangle$ bloquée par $\{E\}$?
- ▶ Chaîne $\langle D, G, C, F, A \rangle$ bloquée par $\{E\}$?

3 Apprentissage de structure

Apprentissage de structure – une 1ère idée

- ▶ **Objectif** : déterminer la structure \mathcal{G} à partir de données \mathbf{D}
- ▶ **Rappel** : \mathcal{G} P-map (perfect map) de P ssi

$$A \perp\!\!\!\perp_P B | C \iff \langle A \perp_{\mathcal{G}} B | C \rangle$$
- ▶ **Algorithme « naïf »** :
 - ▶ créer toutes les structures \mathcal{G} possibles
 - ▶ $\forall \mathcal{G}$, calculer tous les triplets (A, B, C) t.q. $\langle A \perp_{\mathcal{G}} B | C \rangle$
 - ▶ tester si $A \perp\!\!\!\perp_P B | C$ (par exemple, test du χ^2 en utilisant \mathbf{D})
 - ▶ si vrai pour tout triplet (A, B, C) , structure \mathcal{G} trouvée

Problèmes de l'algorithme naïf (1/3)

Théorème – Robinson(1977)

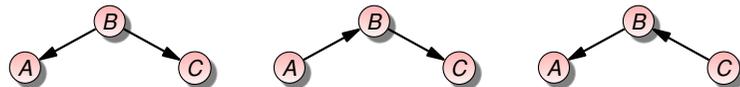
Le nombre de structures \mathcal{G} à n nœuds est super-exponentiel :

$$\#(n) = \begin{cases} 1 & \text{si } n \leq 1 \\ \sum_{i=1}^n (-1)^{i+1} 2^{i(n-1)} C_n^i \times \#(n-1) & \text{si } n > 1 \end{cases}$$

⇒ on ne peut pas tester toutes les structures



3 réseaux bayésiens équivalents (mêmes indep.) :



⇒ ne les compter que pour 1 seul réseau !

⇒ Appliquer l'algorithme dans l'espace des classes d'équivalence de Markov

Classe d'équivalence de Markov

Définition : équivalence de Markov

- ▶ $\mathcal{G}_1, \mathcal{G}_2$ deux structures de RB contenant les mêmes nœuds/variables aléatoires
- ▶ $\mathcal{G}_1, \mathcal{G}_2 \in$ même classe d'équivalence de Markov ssi, pour tous ensembles de variables disjoints A, B, C :

$$\langle A \perp_{\mathcal{G}_1} B | C \rangle \iff \langle A \perp_{\mathcal{G}_2} B | C \rangle$$



Définitions

- ▶ \mathcal{G} structure de RB. **Squelette** de \mathcal{G} obtenu en remplaçant les arcs $X \rightarrow Y$ par des arêtes $X - Y$.
- ▶ $\{\text{nœuds } X, Y, Z\} = \text{v-structure} \iff$ dans \mathcal{G} , $\exists X \rightarrow Y \leftarrow Z$ et $\not\exists X \rightarrow Z$ et $\not\exists Z \rightarrow X$

Théorème – Verma et Pearl (1991)

$\mathcal{G}_1, \mathcal{G}_2 \in$ même classe d'équivalence de Markov ssi même squelette et mêmes v-structures.

Problèmes de l'algorithme naïf (2/3)



Appliquer l'algo dans l'espace des classes d'équivalence

Propriété expérimentale – Gillispie et Perlman (2002)

Ratio $\frac{\text{taille de l'espace des DAG}}{\text{taille de l'espace des classes d'équivalence}} \approx 3,7$

⇒ Pas d'avantage à utiliser les classes d'équivalence

Théorème – Chickering, Heckerman, Meek (2004)

L'apprentissage de structure de RB est NP-hard.

▶ **2 alternatives :**

- ▶ Apprentissage « exact » pour des « petits » RB
- ▶ Apprentissage « approché » ⇒ heuristiques

Problèmes de l'algorithme naïf (3/3)



▶ **2ème problème :** tester si $A \perp\!\!\!\perp_P B | C$ impossible si D petite ou $\Omega_{A \cup B \cup C}$ de grande taille

$A \perp_G B | C \iff (X \perp_G Y | C \forall X \in A, \forall Y \in B)$

▶ Perfect map $\implies A \perp\!\!\!\perp_P B | C \iff (X \perp\!\!\!\perp_P Y | C \forall X \in A, \forall Y \in B)$

Tests d'indépendance

- ▶ Hypothèse (**DAG-faithfulness**) : P représentable par une perfect map \mathcal{G}
- ▶ Ne tester que l'indépendance conditionnelle de couples de variables

Théorème d'Hammersley-clifford

P distribution strictement positive $\implies P$ représentable par une perfect map.



Relations déterministes entre variables $\implies P$ non strictement positive

Apprentissage fondé sur les contraintes

Idée générale de l'apprentissage sous contraintes

- 1 Apprendre le squelette via des tests d'indépendance
- 2a Orienter les v-structures
- 2b Propager ces orientations afin qu'elles ne créent pas de nouvelle v-structure
- 3 Orienter le reste des arêtes sans créer de nouvelle v-structure

▶ Algorithmes classiques :

- ▶ Inductive causation (IC) – Verma et Pearl (1990)
- ▶ PC – Spirtes, Glymour et Scheines (2000)

Bibliographie

- ▶ Chickering D., Heckerman D. et Meek C. (2004) « Large-Sample Learning of Bayesian Networks is NP-Hard », Journal of Machine Learning Research, 5 :1287–1330
- ▶ Geiger D. et Heckerman D. (1997) « A Characterization of the Dirichlet Distribution through Global and Local Parameter Independence », The Annals of Statistics, 25(3) :1344–1369
- ▶ Gillispie S. et Perlman M. (2002) « The size distribution for Markov equivalence classes of acyclic digraph models », Artificial Intelligence, 141 :137–155
- ▶ Robinson, R. (1977) « Counting unlabeled acyclic digraphs », Combinatorial Mathematics V, 622 : :28–43
- ▶ Spirtes E., Glymour C. et Scheines R. (2000) Causation, Prediction and Search, 2nd edition, Springer-Verlag
- ▶ Verma T. et Pearl J. (1990) « Equivalence and synthesis of causal models », Proceedings of UAI, 220–227