

Introduction to graphical models :

Bayesian Networks and Hidden Markov Models Part I

June 2018

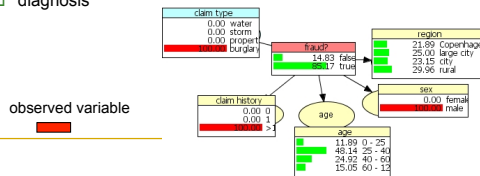
Laurence Likforman-Sulem
IMT/Telecom ParisTech/University Paris-Saclay
likforman@telecom-paristech.fr

Overview

- Part I : Graphical Models
 - Bayesian Networks
 - Dynamic Bayesian Networks (DBN)
 - link with HMMs (Hidden Markov Models)
- Part II : Hidden Markov Models
 - discrete, continuous
 - generative models
 - decoding, training

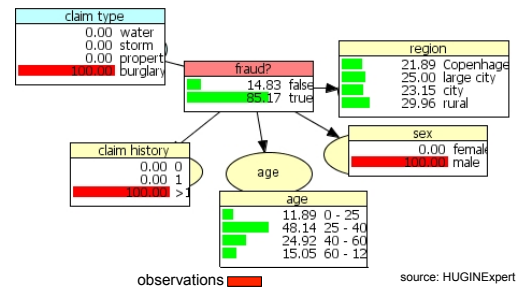
applications

- HMMs and Dynamic Bayesian Networks (DBN)
 - Speech recognition
 - Handwriting recognition
 - Recognition of objects, faces in videos,...
 - Natural Language processing
- Static Bayesian networks
 - Fraud detection (ex: HUGIN)
 - diagnosis



Application: fraud detection

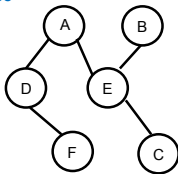
From observation variable(s), infer values of remaining variables



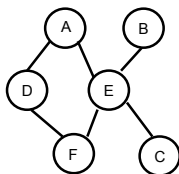
Graphical Models : undirected graphs

- tree : single path between 2 nodes
- multi-connected graph : several paths between nodes

Tree



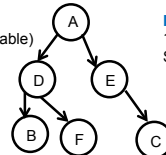
multi-connected graph



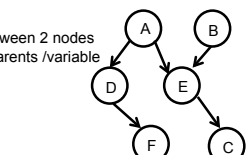
Graphical Models : directed graphs

- Relations : sons, siblings, parents, descendants, ancestors

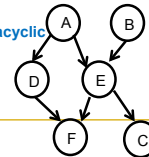
Simple Tree
(1 parent/variable)



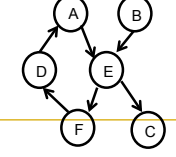
Poly-tree
1 path between 2 nodes
Several parents /variable



multi-connected acyclic Graph

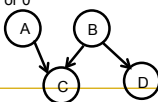


Cyclic Graph

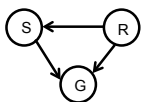


Bayesian network: definition

- Bayesian network: (G, θ)
- $G = (V, E)$ directed acyclic graph (DAG)
 - V : variables (nodes),
 - E : edges: relations between variables (influence of one variable over another)
- θ : parameters (conditional probability tables)
 - Ex: $P(X | \text{parents}(X))$ or $P(X)$ if X has no parent
- Graph : factories joint probability of all variables
 - training: less parameters $1+1+4+2=8$ instead of $2^4-1=15$
 - Inference : computational complexity reduced

$$P(A, B, C, D) = P(A)P(B)P(C|A, B)P(D|B)$$


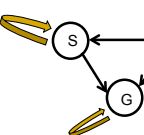
Bayesian network: example

- G: grass is wet (1 ou 0)
 - S: sprinkler on (1 ou 0)
 - R: rain during past night (1 ou 0)
- 
- [from Wikipedia BN]
- $S \rightarrow G$: G is a consequence of S
 - $R \rightarrow S$: R influences whether the sprinkler is on.
(if it has rained, you do not need to open the sprinkler)
 - $R \rightarrow G$: R influences whether the grass is wet
 - Conversely: knowing the value of G «modifies the belief» in S and R : $P(R=1|G=1) > P(R=1)$

Bayesian network: parameters

- Conditional Probability Tables or Distributions (CPT and CPD) : $P(X | \text{parents}(X))$
- G: grass wet (1 or 0)
- S: sprinkler on (1 or 0)
- R: rain during night (1 or 0)

R	$P(S=1 R)$	$P(S=0 R)$
0	0.4	0.6
1	0.01	0.99

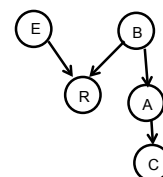


S	R	$P(G=0 S,R)$	$P(G=1 S,R)$
1	1	0.01	0.99
1	0	0.1	0.9
0	1	0.2	0.8
0	0	1.0	0.0

$P(R=0)=0.8$
 $P(R=1)=0.2$

Bayesian network: conditional independence

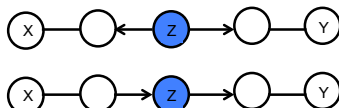
- A node is independent of its non-descendants given its parents
- Graph: binary variables A, B, C, E, R
- C is independent of R, B, E given A



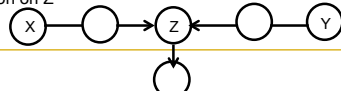
[from Nir Friedman]

conditional independence : d-separation

- A node X is independent from a node Y given a set of observed variables E (E d-separates X and Y)
- If all non directed paths between X and Y are blocked by a variable Z such as :
- Z is in E and there is a chain or divergent connexion on Z

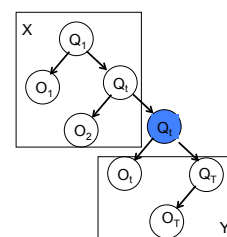


- Z is not observed, nor its descendants, and there is a convergent connexion on Z



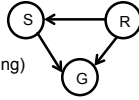
d-separation : set of nodes

- set of nodes X and set of nodes Y are d-separated given observed node Q_t



Bayesian network: inference

- computing $P(\text{variable(s)} \mid \text{observed variable(s)})$
- observed variables : « evidence » or observation
- Inference from observations
- inference algorithms
 - Exact: algorithms trees, polytrees (message passing)
 - stochastic : sampling



Inference in chains: backward variable

- 2 nodes (1 observed) $E = \{Y = y_2\}$
- Parameters : $P(X=x)$ with $\text{dom}\{X\}=\{x_1, x_2\}$, and $P_{Y|X}$ with $\text{dom}\{Y\}=\{y_1, y_2, y_3\}$,

$$P_{Y|X} = \begin{bmatrix} P(Y=y_1|X=x_1) & P(Y=y_2|X=x_1) & P(Y=y_3|X=x_1) \\ P(Y=y_1|X=x_2) & P(Y=y_2|X=x_2) & P(Y=y_3|X=x_2) \end{bmatrix}$$

- backward variable λ :

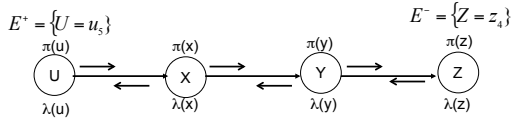
$$\lambda(y) = P(E|Y=y) = [0 \ 0 \ 1 \ 0]^T$$

$$\lambda(x) = P(E|X=x)$$

$$\lambda(x) = P_{Y|X} \cdot \lambda(y) = \begin{bmatrix} P(Y=y_2|X=x_1) \\ P(Y=y_2|X=x_2) \end{bmatrix} \quad \begin{array}{l} \lambda(x) \text{ computed from} \\ \text{« message » } \lambda(y) \text{ sent} \\ \text{by Y to X} \end{array}$$

$$P(X, E) = P(E = e | X = x) P(X = x) = \lambda(x) P(X = x)$$

inference in chains



4 nodes (2 are observed) (Z and U), observations E^+ and E^- :
 $\text{dom}\{Z\}=\{z_1, z_2, z_3, z_4\}$
 $\text{dom}\{U\}=\{u_1, u_2, u_3, u_4, u_5\}$

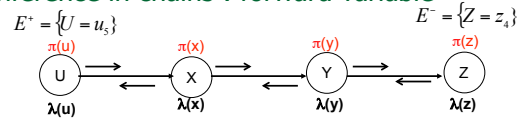
backward variable $\lambda(z) = P(E^-|Z=z) = [0 \ 0 \ 0 \ 1]^T$

$$\lambda(y) = P_{Z|Y} \cdot \lambda(z)$$

$$\text{and } \lambda(x) = P_{Y|X} \cdot \lambda(y)$$

$\lambda(y)$ computed from
 « message » $\lambda(z)$ sent
 to Y by Z
 Then $\lambda(x)$ computed from
 « message » $\lambda(y)$ sent
 to X by Y

inference in chains : forward variable



4 nodes (2 observed)

forward variable $\pi(x) = P(X=x, E^+)$

$$\pi(u) = [0 \ 0 \ 0 \ 0 \ 1]$$

$$\pi(x) = \pi(u) \cdot P_{X|U}$$

$$\pi(y) = \pi(x) \cdot P_{Y|X}$$

$\pi(x)$ computed from
 message $\pi(u)$ sent by U
 to child X

$$P(X=x, E) = \lambda(x) \pi(x)$$

$$P(X=x|E) \propto \lambda(x) \pi(x)$$

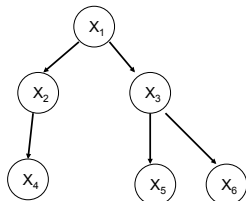
normalisation factor: $1/P(E)$

$$P(E) = \sum_{x \in \text{dom}\{X\}} \lambda(x) \pi(x)$$

inference in trees

- observed variables : $X_1=2, X_4=1, X_5=1, X_6=2$

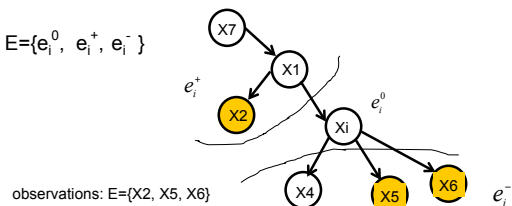
- hidden variables (states) : $X_2, \text{dom}\{X_2\}=\{3,4,5\}, X_3, \text{dom}\{X_3\}=\{6,7\}$



inference in trees

- e_i^- : observed variables in subtrees rooted in X_i 's children
- e_i^0 : observed value of X_i (if X_i observed)
- e_i^+ : all other, observed variables

- $E = \{e_i^0, e_i^+, e_i^-\}$



observations: $E = \{X_2, X_5, X_6\}$

inférence in trees

■ for node X_i (hidden)

□ variable λ (backward) $\lambda(x_i) = P(e_i^0, e_i^- | X_i = x_i)$

□ variable π forward $\pi(x_i) = P(e_i^+, X_i = x_i)$

$$P(E, X_i = x_i) = P(e_i^0, e_i^-, e_i^+, X_i = x_i)$$

$$P(X_i = x_i, E) = \lambda(x_i) \pi(x_i)$$

$$P(E) = \sum_{x_i \in \text{dom}(X_i)} \lambda(x_i) \pi(x_i)$$

$$P(X_i = x_i | E) = \frac{\lambda(x_i) \pi(x_i)}{\sum_{x_i \in \text{dom}(X_i)} \lambda(x_i) \pi(x_i)} \propto \lambda(x_i) \pi(x_i)$$

Message passing inference algorithm

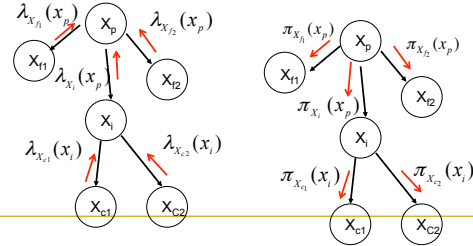
message λ from X_i 's children to X_i : $\lambda_{x_{ci}}(x_i)$

message λ from X_i to single parent X_p : $\lambda_{x_i}(x_p)$

message λ from X_i 's siblings to parent X_p : $\lambda_{x_{j2}}(x_p)$

message π from parent X_p to child X_i : $\pi_{x_i}(x_p)$

message π sent by X_i to its children: $\pi_{x_{ci}}(x_i)$



Inference algorithm : computing λ

□ from leaves to root

■ children X_{ci} send message to parent X_i .

$$\lambda_{x_{ci}}(x_i) = P_{x_{ci}|x_i} \cdot \lambda(x_{ci})$$

■ IF X_i leaf node observed $\lambda(x_i) = [0 \ 1 \ 0]$ (1 at observation position)

■ IF X_i leaf node non observed: $\lambda(x_i) = [1 \ 1 \ 1]$

■ if X_i hidden node (not leaf) : combine children λ messages

$$\lambda(x_i) = \prod_{Y: \text{children}(X_i)} \lambda_{x_{ci}}(x_i) \quad \lambda(x_i) = \lambda_{x_{ci1}}(x_i) \lambda_{x_{ci2}}(x_i)$$

■ si X_i observed node (not leaf) : $\lambda(x_i) = \prod_{Y: \text{children}(X_i)} \lambda_{x_{ci}}(x_i)$, for x_i =corresponding to observation, else 0

Inference algorithm : computing π

□ from root to leaves

□ X_i observed $\pi(x_i) = [0 \ 1 \ 0]$

□ X_i root, non observed $\pi(x_i) = P(X_i = x_i)$

■ else: use message π : $\pi_{x_i}(x_p)$ sent to X_i by unique parent X_p :

$$\pi_{x_i}(x_p) = \pi(x_p) \prod_{Z: \text{siblings of } X_i} \lambda_{x_{j2}}(x_p)$$

■ compute $\pi(x_i) = \pi_{x_i}(x_p) P_{x_i|x_p}$

■ X_i sends messages π to each child

(combining messages λ of siblings of X_i , to parent X_p and π of X_i as a parent)

$$\pi_{x_{ci}}(x_i) = \prod_{X_j: \text{siblings of } X_{ci}} \lambda_{x_j}(x_i) \pi(x_i) = \prod_{X_j} \lambda_{x_j}(x_i) \pi_{x_i}(x_p) P_{x_i|x_p}$$

inference in a DAG

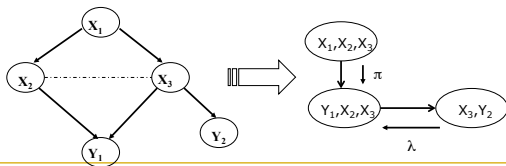
□ convert graph into tree structure (junction tree)

□ inference algorithm : junction tree [Jensen, 96] [Zweig, 2003]

■ moralisation (connect parents)

■ triangulation (cliques construction)

■ connect cliques



training Bayesian networks

■ training with complete data, known structure

- parameter estimation $P(\text{variable} | \text{Parents})$
- maximum likelihood estimation

■ training with incomplete data, known structure

- EM algorithm or gradient descent or stochastic approaches
- MCMC (Gibbs sampling)

■ training with complete data, unknown structure

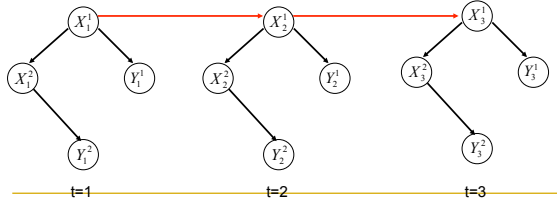
- greedy algorithms

■ training with uncomplete data, unknown structure

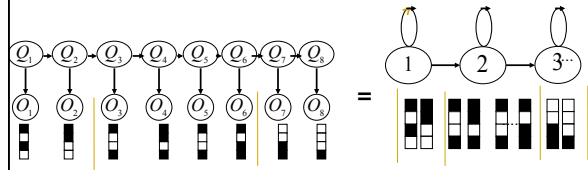
- EM + greedy algorithms

Dynamic Bayesian networks (DBN)

- Extension of static Bayesian networks to dynamic networks
- Stationary processes
 - structure & parameters do not vary through time
 - several state variables and observation variables at each time step.



HMM= special case of DBN



- HMM: Hidden Markov Model
- DBN: tree
- 1 state variable + 1 observation variable at each time step t

$(Q_t)_{1 \leq t \leq T}$: state variable (hidden)

$(O_t)_{1 \leq t \leq T}$: observation variable generated by state variable

références

- Wikipedia BN: http://en.wikipedia.org/wiki/Bayesian_network
- A. W. Moore, *Bayes Nets for representing and reasoning about uncertainty*, 2001
 - <http://www.cs.cmu.edu/~7Eawm/tutorials>
- S. Davis, A. Moore Bayesian networks: independencies and inference
 - <http://www.cs.cmu.edu/~awm/tutorials>
- K. Murphy, BayesNet Toolbox for matlab <https://code.google.com/p/bnt/>
- N. Friedman, D. Koller, Learning Bayesian Networks from data
- G. Zweig, Speech Recognition with Dynamic Bayesian Networks, Phd thesis, 1998, Univ. of California, Berkeley.
- P. Leray, Réseaux Bayésiens-Apprentissage de la structure
 - <http://asi.unsa-roven.fr/enseignants/~pleray/RB2003/2-ApprentissageStructure.pdf>
- P. Naim, P.-H. Willemin, P. Leray, O. Pourret, A. Becker, Les réseaux Bayésiens, Eyrolles, 2007.
- M. Sigelle, Bases de la Reconnaissance des Formes: Chaînes de Markov et Modèles de Markov Cachés, chapitre 7, Polycopié Telecom ParisTech, 2012.
- J. Pearl, Probabilistic Reasoning in intelligent systems: networks of plausible inference networks, 1988.
- L. Likforman-Sulem, E. Barney Smith, Reconnaissance des Formes: théorie et pratique sous matlab, Ellipses, TechnoSup, 2013.