

# Bayesian Networks Learning

Fabrizio Riguzzi  
Engineering Department



## Summary

---

- Probability theory
- Conditional independence
- Definition of Bayesian network
- Inference
- Learning
- Logic and probability

2

## Uncertainty

---

- Reasoning requires simplifications:
  - Birds fly
  - Smoke suggests fire
- Treatment of exceptions
- How to reason from uncertain knowledge?

3

## How to Perform Inference?

---

- Use non-numerical techniques
  - Logician: non monotonic logic
- Assign to each proposition a numerical measure of uncertainty
  - Neo-probabilist: use probability theory
  - Neo-calculist: use other theories:
    - fuzzy logic
    - certainty factors
    - Dempster-Shafer

4

## Probability Theory

---

- A: Proposition,
  - Ex: A=The coin will land heads
- P(A): probability of A
- Frequentist approach: probability as relative frequency
  - Repeated random experiments
  - P(A) is the fraction of experiments in which A is true
- Bayesian approach: probability as a degree of belief
- Example: B=burglary tonight

5

## Axioms of Probability Theory

---

$$0 \leq P(A) \leq 1$$

$$P(\text{Sure Proposition}) = 1$$

$$P(A \vee B) = P(A) + P(B)$$

if A and B are mutually exclusive

6

## Probability Rules

---

- Any event A can be written as the or of two disjoint events (A and B) and (A and  $\neg B$ )

$$P(A) = P(A, B) + P(A, \neg B) \quad \text{marginalization/sum rule}$$

- Where  $P(A, B) = P(A \wedge B)$  is called the **joint probability** of A and B
- In general, if  $B_i, i=1, 2, \dots, n$  is a set of exhaustive and mutually exclusive propositions

$$P(A) = \sum_i P(A, B_i)$$

- Moreover

$$P(A) + P(\neg A) = 1$$

7

## Conditional Probabilities

---

- $P(A|B)$  = belief of A given that I know B
- Relation to  $P(A, B)$

$$P(A, B) = P(A|B)P(B) \quad \text{product rule}$$

$$P(A|B) = \frac{P(A, B)}{P(B)}$$

8

## Bayes Theorem

---

- Relationship between  $P(A|B)$  and  $P(B|A)$

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

- $P(A)$ : **prior probability**
- $P(A|B)$ : **posterior probability** (after learning B)

9

## Conditional Independence

---

- If  $P(A|B)=P(A)$  we say that A and B are independent
- If  $P(A|B,C)=P(A|C)$  we say that A and B are conditionally independent given C

10

## Chain Rule

---

- n events  $E_1, \dots, E_n$
- Joint event  $(E_1, \dots, E_n)$

$$P(E_1, \dots, E_n) = P(E_n | E_{n-1}, \dots, E_1) P(E_{n-1}, \dots, E_1)$$

$$P(E_1, \dots, E_{n-1}) = P(E_{n-1} | E_{n-2}, \dots, E_1) P(E_{n-2}, \dots, E_1)$$

...

- Chain rule:

$$P(E_1, \dots, E_n) = P(E_n | E_{n-1}, \dots, E_1) \dots P(E_2 | E_1) P(E_1) =$$

$$\prod_{i=1}^n P(E_i | E_{i-1}, \dots, E_1)$$

11

## Multivalued Hypothesis

---

- Propositions can be seen as binary variables, i.e. variables taking values true or false
  - Burglary B: true or false
- Generalization: multivalued variables
  - Semaphore S, values: green, yellow, red
  - Propositions are a special case with two values

12

## Discrete Random Variables

---

- Variable  $V$ , values  $v_i$   $i=1, \dots, n$
- $V$  is also called a **discrete random variable**
- $V=v_i$  is a proposition
- Propositions  $V=v_i$   $i=1, \dots, n$  exhaustive and mutually exclusive
- $P(v_i)$  stands for  $P(V=v_i)$
- $V$  is described by the set  $\{P(v_i)|i=1, \dots, n\}$ , the **probability distribution** of  $V$ , indicated with  $P(V)$

13

## Notation

---

- We indicate with  $v$  a generic value of  $V$
- Set or vector of variables:  $\mathbf{V}$ , values  $\mathbf{v}$

14

## Marginalization

---

- Multivalued variables  $A$  and  $B$
- $b_i$   $i=1, \dots, n$  values of  $B$

$$P(a) = \sum_i P(a, b_i)$$

- Or

$$P(a) = \sum_b P(a, b)$$

- In general

$$P(x) = \sum_y P(x, y) \quad \text{sum rule}$$

15

## Conditional Probabilities

---

- $P(a|b)$  = belief of  $A=a$  given that know  $B=b$
- Relation to  $P(a, b)$

$$P(a, b) = P(a|b) P(b) \quad \text{product rule}$$

$$P(a|b) = \frac{P(a, b)}{P(b)}$$

- Bayes theorem

$$P(a|b) = \frac{P(b|a) p(a)}{P(b)}$$

16

## Continuous Random Variables

---

- A multivalued variable  $V$  that takes values from a real interval  $[a,b]$  is called a **continuous random variable**
- $P(V=v)=0$ , we want to compute  $P(c \leq V \leq d)$
- $V$  is described by a **probability density function**  $\rho: [a,b] \rightarrow [0,1]$
- $\rho(v)$  is such that

$$P(c \leq V \leq d) = \int_c^d \rho(v) dv$$

17

## Properties of Continuous Random Variables

---

- The same as those of discrete random variables where summation is replaced by integration:

- Marginalization (sum rule)

$$\rho(x) = \int \rho(x, y) dy$$

- Conditional probability (product rule)

$$\rho(x, y) = \rho(x|y)\rho(y)$$

....

18

## Mixed Distribution

---

- We can have a conjunction of discrete and continuous variables
- Joint: if one of the variables is continuous, the joint is a density:
  - X discrete, Y continuous:  $\rho(x,y)$
- Conditional joint:
  - X discrete, Y continuous:  $P(x|y)$
  - X discrete, Y continuous, Z discrete:  $\rho(x,y|z)$

19

## Domain Modeling

---

- We use a set of random variables to describe the domain of interest
- Example: home intrusion detection system, variables:
  - Earthquake E, values  $e_1=no$ ,  $e_2=moderate$ ,  $e_3=severe$
  - Burglary B, values:  $b_1=no$ ,  $b_2=yes$  through door,  $b_3=yes$  through window
  - Alarm A, values  $a_1=no$ ,  $a_2=yes$
  - Neighbor call N, values  $n_1=no$ ,  $n_2=yes$

20

## Inference

---

- We would like to answer the following questions
  - What is the probability of a burglary through the door? (compute  $P(b_2)$ , belief computation)
  - What is the probability of a burglary through the window given that the neighbor called ? (compute  $P(b_2|n_2)$ , belief updating)

21

## Inference

---

- What is the probability of a burglary through the door given that there was a moderate earthquake and the neighbor called ? (compute  $P(b_2|n_2,e_2)$ , belief updating )
- What is the probability of a burglary through the door and of the alarm ringing given that there was a moderate earthquake and the neighbor called ? (compute  $P(a_2,b_2|n_2,e_2)$ , belief updating)
- What is the most likely value for burglary given that the neighbor called ( $\text{argmax}_b P(b|n_2)$ , belief revision)

22

## Types of Problems

---

- Diagnosis:  $P(\text{cause}|\text{symptom})=?$
- Prediction:  $P(\text{symptom}|\text{cause})=?$
- Classification:  $\text{argmax}_{\text{class}} P(\text{class}|\text{data})$

23

## Inference

---

- In general, we want to compute the probability  $P(\mathbf{q}|\mathbf{e})$ 
  - of a query  $\mathbf{q}$  (assignment of values to a set of variables  $\mathbf{Q}$ )
  - given the evidence  $\mathbf{e}$  (assignment of values to a set of variables  $\mathbf{E}$ )

24

## Joint Probability Distribution

---

- The **joint probability distribution** (jpd) of a set of variables  $\mathbf{U}$  is given by  $P(\mathbf{u})$  for all values  $\mathbf{u}$
- For our example
  - $\mathbf{U}=\{E,B,A,N\}$
  - We have the jpd if we know  $P(\mathbf{u})=P(e,b,a,n)$  for all the possible values  $e, b, a, n$ .

25

## Inference

---

- If we know the jpd, we can answer all the possible queries:

$$P(\mathbf{q}|\mathbf{e}) = \frac{P(\mathbf{q}, \mathbf{e})}{P(\mathbf{e})} = \frac{\sum_{x, X \in \mathbf{U} \setminus \mathbf{Q} \setminus \mathbf{E}} P(x, \mathbf{q}, \mathbf{e})}{\sum_{x, X \in \mathbf{U} \setminus \mathbf{E}} P(x, \mathbf{e})}$$

26

## Problem

---

- If we have  $n$  binary variables ( $|\mathbf{U}|=n$ ), knowing the jpd requires storing  $O(2^n)$  different values.
- Even if had the space to store all the  $2^n$  different values, computing  $P(\mathbf{q}|\mathbf{e})$  would require  $O(2^n)$  operations
- Impractical for real world domains
- How to avoid the space and time problems? Use conditional independence assertions

27

## Conditional Independence

---

- $\mathbf{X}, \mathbf{Y}, \mathbf{Z}$  vectors of multivalued variables
- $\mathbf{X}$  and  $\mathbf{Y}$  are **conditionally independent** given  $\mathbf{Z}$  if

$$P(\mathbf{x}|\mathbf{y}, \mathbf{z}) = P(\mathbf{x}|\mathbf{z}) \text{ whenever } P(\mathbf{y}, \mathbf{z}) > 0$$

- We write  $I\langle \mathbf{X}, \mathbf{Z}, \mathbf{Y} \rangle$
- Special case:  $\mathbf{X}$  and  $\mathbf{Y}$  are **independent** if

$$P(\mathbf{x}|\mathbf{y}) = P(\mathbf{x}) \text{ whenever } P(\mathbf{y}) > 0$$

28

## Chain Rule

---

- n random variables  $X_1, \dots, X_n$
- Let  $\mathbf{U} = \{X_1, \dots, X_n\}$
- Joint event  $\mathbf{u} = (x_1, \dots, x_n)$
- Chain rule:

$$\begin{aligned} P(\mathbf{u}) &= P(x_1, \dots, x_n) \\ &= P(x_n | x_{n-1}, \dots, x_1) \dots P(x_2 | x_1) P(x_1) \\ &= \prod_{i=1}^n P(x_i | x_{i-1}, \dots, x_1) \end{aligned}$$

29

## Conditional Independence

---

- $\Pi_i$  is a subset of  $\{X_{i-1}, \dots, X_1\}$  such that
- $X_i$  is conditionally independent of  $\{X_{i-1}, \dots, X_1\} \setminus \Pi_i$  given  $\Pi_i$

$$P(x_i | x_{i-1}, \dots, x_1) = P(x_i | \pi_i)$$

- where  $\pi_i$  is a set of values for  $\Pi_i$
- $\Pi_i$  parents of  $X_i$

30

## Conditional Independence

---

- Knowing  $\Pi_i$  for all i we could write

$$\begin{aligned} P(\mathbf{u}) &= P(x_1, \dots, x_n) \\ &= P(x_n | x_{n-1}, \dots, x_1) \dots P(x_2 | x_1) P(x_1) \\ &= P(x_n | \pi_n) \dots P(x_2 | \pi_2) P(x_1 | \pi_1) \\ &= \prod_{i=1}^n P(x_i | \pi_i) \end{aligned}$$

31

## Conditional Independence

---

- In order to compute  $P(\mathbf{u})$  we have to store

$$P(x_i | \pi_i)$$

- for all values  $x_i$  and  $\pi_i$
- $P(x_i | \pi_i)$ : Conditional probability table
- If  $\Pi_i$  is much smaller than the set  $\{X_{i-1}, \dots, X_1\}$ , then we have huge savings
- If k is the maximum number of parents of a variable, then storage is  $O(n2^k)$  instead of  $O(2^n)$

32



## Graphical Representation

- We can represent the conditional independence assertions using a directed graph network with a node per variable
- $\Pi_i$  is the set of parents of  $X_i$
- The graph is acyclic

33

## Example Network

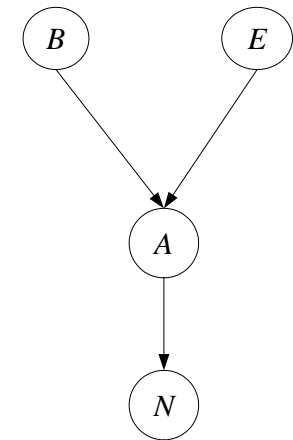
- Variable order: E,B,A,N
- Independences

$$P(e)$$

$$P(b|e) = P(b)$$

$$P(a|b, e) = P(a|b, e)$$

$$P(n|a, b, e) = P(n|a)$$



34

## Bayesian Network

- A Bayesian network [Pearl 85] (BN)  $B$  is a couple  $(G, \Theta)$  where
  - $G$  is a directed acyclic graph (DAG)  $(V, E)$  where
    - $V$  is a set of vertices  $\{X_1, \dots, X_n\}$
    - $E$  is a set of edges, i.e. A set of couples  $(X_i, X_j)$
    - $\langle X_1, \dots, X_n \rangle$  is a topological sort of  $G$ , i.e.  $(X_i, X_j) \in E \Rightarrow i < j$
  - $\Theta$  is a set of conditional probability tables (cpt)
 
$$\{\theta_{x_i|\pi_i} | i=1, \dots, n, x_i \in X_i, \pi_i \in \Pi_i\}$$
  - where  $\Pi_i$  is the set of parents of  $X_i$

35

## Bayesian Network

- A BN  $(G, \Theta)$  **represents** a jpd  $P$  iff
  - each variable is independent of its predecessors given its parents in  $G$ 

$$P(x_i | x_{i-1}, \dots, x_1) = P(x_i | \pi_i)$$
  - $\theta_{x_i|\pi_i} = P(x_i | \pi_i)$  for all  $i$  and  $\pi_i$
- In this case

$$P(x_1, \dots, x_n) = \prod_{i=1}^n P(x_i | \pi_i)$$

$$= \prod_{i=1}^n \theta_{x_i|\pi_i}$$

36

## How to Build a Bayesian Network

---

- Choose an ordering  $X_1 \dots X_n$  for the variables
- For  $i = 1$  to  $n$ :
  - Add  $X_i$  node to the network
  - Set  $\Pi_i$  to be a minimal subset of  $\{X_1 \dots X_{i-1}\}$  such that we have conditional independence of  $X_i$  and all other members of  $\{X_1 \dots X_{i-1}\}$  given  $\Pi_i$
  - Assign a value to  $P(x_i | \pi_i)$  for all the values of  $x_i$  and  $\pi_i$

37

## Building a Bayesian Network

---

- Usually the expert consider a variable  $X$  as a child of  $Y$  if  $Y$  is a **direct cause** of  $X$
- Correlation and causality are related but are **not** the same thing
  - See the book [Pearl 00]

38

## Pathfinder system [Suermondt et al. 90]

---

- Diagnostic system for lymph-node diseases.
- 60 diseases and 100 symptoms and test-results.
- 14,000 probabilities
- Expert consulted to make net.
- 8 hours to determine variables.
- 35 hours for net topology.
- 40 hours for probability table values.

39

## Pathfinder system [Suermondt et al. 90]

---

- Pathfinder is now outperforming the world experts in diagnosis.
- Being extended to several dozen other medical domains.

40

## Inference with Bayesian Networks

- With a Bayesian Network we save space, do we also save time?
- Do we have to use the formula

$$P(\mathbf{q}|\mathbf{e}) = \frac{\sum_{x, X \in V \setminus Q \setminus E} P(x, \mathbf{q}, \mathbf{e})}{\sum_{x, X \in V \setminus E} P(x, \mathbf{e})}$$

- to compute  $P(\mathbf{q}|\mathbf{e})$ ?

41

## Inference with Bayesian Networks

- There are quicker algorithms
  - Exact methods for polytrees
    - Belief propagation
  - Exact methods for general networks
    - Junction tree
    - Variable elimination
  - Approximate methods for general networks:
    - Stochastic sampling
    - Loopy belief propagation
    - Variational methods,

42

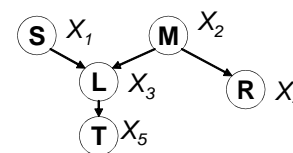
## Complexity of Inference

- Exact inference with BN is #P-complete
- #P-complete: a special case of NP-complete problems
  - The answer to a #P-complete problem is the number of solutions to a NP-complete problem

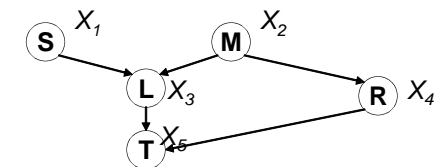
43

## Polytrees

A polytree is a directed acyclic graph in which no two nodes have more than one path between them.



A polytree



Not a polytree

- i.e. There are no cycles in the corresponding undirected graph

44

## Belief Propagation [Pearl 88]

- To compute  $P(x|e)$  write

$$P(x|e) = \alpha \lambda(x) \pi(x)$$

- where  $\alpha$  is a normalizing constant and
  - $\pi(x)$  represents the support to the assertion  $X=x$  by the non-descendants of  $X$
  - $\lambda(x)$  represents the support to the assertion  $X=x$  by the descendants of  $X$

45

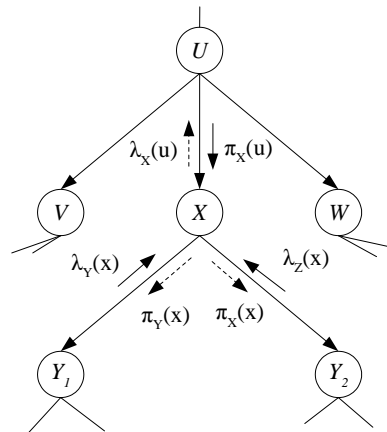
## Belief Propagation

- Nodes exchange messages with their neighbors
- $\pi(x)$  and  $\lambda(x)$  are computed from message received respectively from the parents and the children of  $X$
- When a node is activated:
  - It reads the incoming messages
  - It updates  $\pi(x)$  and  $\lambda(x)$
  - It updates  $P(x|e)$
  - It generates the new messages to be sent to their parents and children

46

## Messages Received

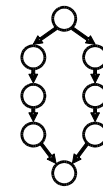
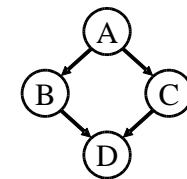
- Node  $X$ ,
- Parents  $U_i$
- Children  $Y_j$



47

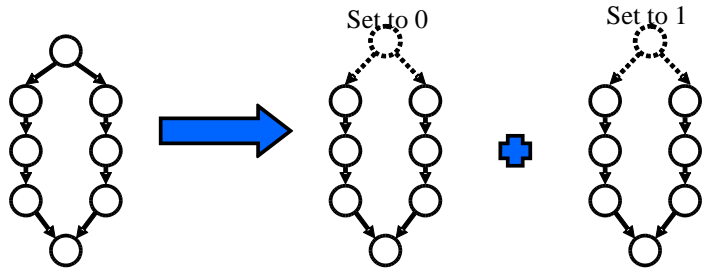
## General Networks

- Networks that have a cycle in their undirected version
- Two possibilities
  - Conditioning
  - Clustering



48

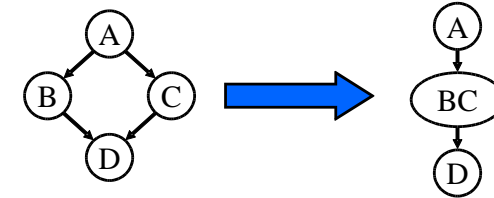
## Conditioning



49

## Clustering

- Group together variables so that the resulting network is a polytree and use belief propagation



- Problem: how to find a good clustering?

50

## Join Trees

- Technique for clustering variables
- Steps:
  - Obtain an undirected version of the network
  - Perform a graph operation on it (triangulation)
  - Each clique is a compound variable
  - Add direction to the edges

51

## Junction Tree

- The resulting inference algorithm [Lauritzen, Spiegelhalter 1988] is called
  - Junction tree algorithm (jt), or
  - Clique propagation

52

## Approximate Methods

---

- Sampling:
  - Generate N samples from BN
  - Count:  $N_e$ : samples that satisfy  $\mathbf{e}$ ,  $N_{q_e}$  samples that satisfy  $\mathbf{q}, \mathbf{e}$
  - $P(\mathbf{q}|\mathbf{e}) = N_{q_e} / N_e$
- Loopy belief propagation:
  - bp in networks with cycles
  - Experiments have shown that it converges to good quality solutions

53

## Sampling

---

- Let  $X_1, \dots, X_n$  be a topological sort of the variables
- For  $i=1$  to  $n$ 
  - Find parents, if any, of  $X_i$ . Call them  $X_{p(i,1)}, X_{p(i,2)}, \dots, X_{p(i,p(i))}$ .
  - Recall the values that those parents were randomly given:  $x_{p(i,1)}, x_{p(i,2)}, \dots, x_{p(i,p(i))}$ .
  - Look up in the cpt for:  
 $P(X_i=x_i \mid X_{p(i,1)}=x_{p(i,1)}, X_{p(i,2)}=x_{p(i,2)} \dots X_{p(i,p(i))}=x_{p(i,p(i))})$
  - Randomly choose  $x_i$  according to this probability

54

## Problems in Building BN

---

- Assessing conditional independence is not always easy for humans
- Usually done on the basis of causal information
- Assigning a number to each cpt entry is also difficult for humans

55

## Problems in Building BN

---

- Often we do not have an expert but we are given a set of observations  $D = \{\mathbf{u}^1, \dots, \mathbf{u}^N\}$
- $\mathbf{u}^j$  is an assignment to all the variables  $\mathbf{U} = \{X_1, \dots, X_n\}$
- How to infer the parameters and/or the structure from  $D$ ?

56

## Learning

---

- We want to find a BN over  $\mathbf{U}$  such that the probability of the data  $P(D)$  is maximized
- $P(D)$  is also called the **likelihood** of the data
- We assume that all the samples are **independent and identically distributed** (iid) so

$$P(D) = \prod_i^N P(\mathbf{u}^i)$$

- Often the natural log of  $P(D)$  (**log likelihood**) is considered

$$\log P(D) = \sum_i^N \log P(\mathbf{u}^i)$$

57

## Learning BN

---

- Tasks
  - Computing the parameters given a fixed structure or
  - finding the structure and the parameters
- Properties of data:
  - complete data: in each data vectors  $\mathbf{u}^i$ , the values of all the variables are observed
  - incomplete data

58

## Parameter Learning from Complete Data

---

- Parameters to be learned

$$\theta_{x_i|\pi_i} = P(x_i|\pi_i)$$

- for all  $x_i, \pi_i, i=1, \dots, n$
- The values of the parameters that maximize the likelihood can be computed in closed form

59

## Maximum Likelihood Parameters

---

- Given by relative frequency
- If  $N_y$  be the number of vectors of  $D$  where  $\mathbf{Y}=\mathbf{y}$ .

$$\theta_{x_i|\pi_i} = \frac{N_{x_i, \pi_i}}{N_{\pi_i}}$$

- Counting: for each  $i$ , for each value  $\pi_i$  we must collect

$$C_{\pi_i} = \langle N_{x_i^1, \pi_i}, \dots, N_{x_i^{v(i)}, \pi_i} \rangle$$

- where  $v(i)$  is the number of values of  $X_i$

60

## Structure Learning from Complete Data

- Perform a local search in the space of possible structures
- HGC algorithm [Heckerman, Geiger, Chickering 95]:
  - Start with a structure BestG' (possibly empty)
  - Repeat
    - BestG=BestG'
    - Let Ref={G'|G' is obtained from BestG' by adding, deleting or reversing an arc}
    - Let BestG'=argmax<sub>G'</sub> {score(G')|G' ∈ Ref}
  - while score(BestG')-score(BestG)>0

61

## Structure Score

$$\text{score}(G) = P(D|G)$$

$$\begin{aligned} P(D|G) &= \int \rho(D, \Theta|G) d\Theta \\ &= \int P(D|\Theta, G) \rho(\Theta) d\Theta \end{aligned}$$

- where

$$\begin{aligned} \rho(\Theta) &= \prod_{i, \pi_i} \rho(\theta_{\pi_i}) \\ \theta_{\pi_i} &= \langle \theta_{x_i^1 | \pi_i}, \dots, \theta_{x_i^{v(i)} | \pi_i} \rangle \end{aligned}$$

- and  $\rho(\theta_{\pi_i})$  is the prior density of the vector  $\theta_{\pi_i}$

62

## Prior Density of the Parameters

- A common choice for the form of the prior density is the **Dirichlet probability density**
- In this case  $\rho(\theta_{\pi_i})$  is described by  $v(i)$  parameters

$$C'_{\pi_i} = \langle N'_{x_i^1, \pi_i}, \dots, N'_{x_i^{v(i)}, \pi_i} \rangle$$

- Prior counts: it is as if we had previously observed some data on which the counts are  $N'_{x_i, \pi_i}$

63

## Structure Score

- If the priors for the parameters are Dirichlet, then the score is called BD (for Bayesian Dirichlet) and

$$BD(G) = \sum_i BD_i(G)$$

- where  $BD_i(G)$  depends only on  $C_i$  and  $C'_i$ , the counts for the family of  $X_i$

$$\begin{aligned} C_i &= \langle C_{\pi_i^1}, \dots, C_{\pi_i^{v(i)}} \rangle \\ C'_i &= \langle C'_{\pi_i^1}, \dots, C'_{\pi_i^{v(i)}} \rangle \end{aligned}$$

64



## Structure Score

---

- $BD(G)$  is **decomposable**:
  - It can be computed independently for each family
- Each edge operation involves
  - 1 family (addition, deletion) or
  - 2 families (reversal)
- $BD(G')$  can be quickly computed from  $BD(\text{Best}G)$  by changing only the score of the affected families

65

## Parameter Learning from Incomplete Data

---

- The maximum likelihood parameters cannot be computed in closed form
- An iterative algorithm is necessary: the EM algorithm
- Finds a (possibly) local maximum of the likelihood

66

## EM Algorithm

---

- Initialize the parameters at random  $\Theta$
- Repeat
  - Expectation step:
    - compute the probability of each value of the missing attributes using  $(G, \Theta)$  and inference
    - Obtain a new dataset  $D'$  by completing  $D$  according to the probabilities computed above
  - Compute  $\Theta$  by maximum likelihood on  $D'$ 
    - Relative frequency

67

## Structure Learning from Incomplete Data

---

- There is no decomposable score
- HGC would not be efficient
- Structural EM:
  - Start with a structure  $\text{Best}G'$  (possibly empty)
  - Repeat
    - $\text{Best}G = \text{Best}G'$
    - Compute the parameters of  $\text{Best}G$  with EM
    - Optimize a lower bound of the likelihood of the observed data
    - Let  $\text{Best}G'$  the optimum
  - Until no improvement

68

## Applications of BN

- Monitoring of emergency care patients
- Model of barley crops yield.
- Diagnosis of carpal tunnel syndrome
- Insulin dose adjustment (DBN) in diabetes
- Predicting hails in northern Colorado.
- Evaluating insurance applications

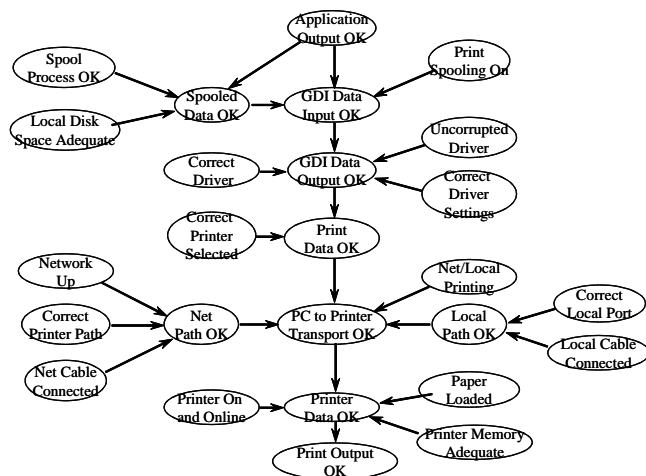
69

## Applications of BN

- Deciding on the amount of fungicides to be used against attack of mildew in wheat.
- Assisting experts of electromyography.
- Pedigree of breeding pigs.
- Modeling the biological processes of a water purification plant.
- Printer troubleshooting (Microsoft Windows)

70

## Printer Troubleshooting (Windows 95)



71

## Applications

- Office Assistant in MS Office (“smiley face”)
  - Bayesian network based free-text help facility
  - help based on past experience (keyboard/mouse use) and task user is doing currently

72

## Markov Networks (MN)

---

- Approach alternative to BN
- Undirected graph
- Conditional independence represented by graph separation
- Probability distribution as the product of a set of **potentials** (functions of a subset of variables) divided by a normalization constant
- One potential per clique

73

## Markov Network

---

- Inference:
  - Algorithms similar to those for BN (bp, ct, ve, ss..)
  - Same complexity
- MN can represent some independences that BN can not represent and vice versa
- Advantage: we do not have to avoid cycles
- Disadvantage: MN parameters are more difficult to interpret

74

## Combination of Logic and Probability

---

- BN are not able to deal with domains containing multiple entities connected by complex relationships
- Logic is not able to represent uncertainty efficiently
- Combination: active research area with many different proposals
- The most common approach is to design a new language and then provide a translation into BN or MN for defining the semantics, performing inference and learning

75

## Some Logical-Probabilistic Languages

---

- Probabilistic Relational Models  $\rightarrow$ BN
- Markov Logic Network  $\rightarrow$ MN
- Bayesian Logic Programs  $\rightarrow$ BN
- Logic Programs with Annotated Disjunctions  $\rightarrow$ BN
- Relational Markov Networks  $\rightarrow$ MN
- CLP(BN)  $\rightarrow$ BN

76

## CLP(BN) [Costa et al 03]

- Based on Prolog
- Variables in a CLP(BN) program can be random
- Their values, parents and CPTs are defined with the program
- To answer a query with uninstantiated random variables, CLP(BN) builds a BN and performs inference
- The answer will be a probability distribution for the variables

77

## Example

```
.....
registration_grade(Key, Grade):-
registration(Key, CKey, SKey),
course_difficulty(CKey, Dif),
student_intelligence(SKey, Int),
{ Grade = grade(Key) with
  p([a,b,c,d],
  %h h h m h l m h m m m l l h l m l l
  [0.20,0.70,0.85,0.10,0.20,0.50,0.01,0.05,0.10,
   0.60,0.25,0.12,0.30,0.60,0.35,0.04,0.15,0.40,
   0.15,0.04,0.02,0.40,0.15,0.12,0.50,0.60,0.40,
   0.05,0.01,0.01,0.20,0.05,0.03,0.45,0.20,0.10 ],
  [Int,Dif]))
}.
.....
```

78

## Inference

```
?- [school_32].
?- registration_grade(r0,G).
p(G=a)=0.4115,
p(G=b)=0.356,
p(G=c)=0.16575,
p(G=d)=0.06675 ?
?- registration_grade(r0,G),
   student_intelligence(s0,h).
p(G=a)=0.6125,
p(G=b)=0.305,
p(G=c)=0.0625,
p(G=d)=0.02 ?
```

79

## Availability

- CLP(BN) is included in Yap prolog
- <http://www.dcc.fc.up.pt/~vsc/Yap/>
- It can use either junction tree or variable elimination for inference

80

## Logic Programs with Annotated Disjunction

---

- [Vennekens et al. 04]
- Minimal extension of logic programming to allow the representation of uncertainty
- Clauses of the form

$$h_1:\alpha_1 ; \dots ; h_n:\alpha_n :- b_1, \dots, b_m$$

- where  $h_i$  are atoms,  $b_i$  are literals and  $\alpha_i$  are probabilities such that

$$\sum_{i=1}^n \alpha_i \leq 1$$

81

## Semantics

---

- Each clause can be seen as an experiment: if  $b_1, \dots, b_m$  is true then  $h_i$  is true with probability  $\alpha_i$  or no  $h_i$  is true with probability  $1 - \sum_i \alpha_i$
- Each ground atom is seen as random variable with values true and false
- We want to assign probabilities to queries (conjunctions of ground atoms), possibly conditioned on some evidence

82

## Semantics

---

- Given an LPAD  $T$ , generate its grounding  $T'$
- An instance of  $T$  is a normal logic program obtained by selecting one head from each clause of  $T'$
- The probability of an instance is obtained by multiplying the probability of each head selected
- The probability of a query  $Q$  is given by the sum of the probabilities of the instances that have  $Q$  as consequence

83

## Example

---

```
heads(Coin):0.5 ; tails(Coin):0.5 :-  
  toss(Coin), \+ biased(Coin).
```

```
heads(Coin):0.6 ; tails(Coin):0.4 :-  
  toss(Coin), biased(Coin).
```

```
biased(Coin):0.1 ; fair(Coin):0.9.
```

```
toss(coin).
```

```
P(heads(coin))=0.51
```

```
P(heads(coin)|biased(coin))=0.6
```

84

## Conversion to Bayesian Networks

---

- An LPAD can be converted to a BN that has
  - One boolean variable per ground atom
  - One variable  $ch_r$  per ground clause  $r$ , with the ground atoms in the head plus null as values
- The dependencies are defined as follows:
  - Ground atom  $a$  depends on all the clause variables that have  $a$  in the head
  - The CPT assign probability 1 to  $a$  if at least one parent is equal to  $a$  and 0 otherwise

85

## Conversion to Bayesian Networks

---

- $ch_r$  depends on the variables that appear in the body of  $r$
- CPT:
  - $P(ch_r=h_i)=\alpha_i$ ,  $P(ch_r=null)=1-\sum_i \alpha_i$  if the body is true
  - $P(ch_r=null)=1$  if the body is false

86

## Inference with LPADs

---

- Convert to BN and use BN inference
  - Problem: the grounding may be very large
- Compute all the possible derivations and compute the probability that one of these derivations is possible [Riguzzi 07]
- Suite of reasoning tools for LPADs: cplint  
<http://www.ing.unife.it/software/cplint/>
- It is included in the CVS version of Yap

87

## Learning LPADs

---

- Data  $D$ : set of interpretations (i.e. sets of ground atoms),
- Task: find the parameters of an LPAD that maximize the likelihood of  $D$ :
- Task: find the parameters and the structure of an LPAD that maximize the likelihood of  $D$

88

## Learning Parameters

---

- ME-compliant LPAD: every couple of ground clauses that share an atom in the head have mutually exclusive bodies
- If an LPAD is ME-compliant then the parameters can be computed in closed form as a ratio of counts [Riguzzi 04]

$$\alpha_i = P(h_i | \text{body})$$

- Otherwise [Blookey, Meert 06]
  - Convert the LPAD to a BN
  - Use EM since the  $ch_r$  variables are unobserved in D

89

## Learning the Structure

---

- If the LPAD is ME-compliant then the structure can be learned by solving a mixed integer programming problem
  - ALLPAD system [Riguzzi 08]
- Otherwise [Blookey, Meert, 07]
  - Use Structural EM to learn a BN
  - Convert to LPAD

90

## BN Software

---

- List of BN software  
<http://www.cs.ubc.ca/~murphyk/Software/bnsoft.html>
- BNT: inference and learning, Matlab, open source
- MSBNx: inference, by Microsoft, free closed source
- OpenBayes: inference and learning, Python, open source
- BNJ: inference and learning, Java, open source
- Weka: learning, Java, open source

91

## Resources

---

- Probabilistic Reasoning in Intelligent Systems by Judea Pearl. Morgan Kaufmann: 1998.
- Probabilistic Reasoning in Expert Systems by Richard Neapolitan. Wiley: 1990.
- List of BN Models and Datasets  
<http://www.cs.huji.ac.il/labs/compbio/Repository/>

92

## Acknowledgments

---

- Some slides from
  - Andrew Moore's tutorials  
<http://www.autonlab.org/tutorials/>
  - Irina Rish and Moninder Singh's tutorial  
<http://www.research.ibm.com/people/r/rish/>

93

## References

---

- [Pearl 85] Pearl, J., "Bayesian Networks: a Model of Self-Activated Memory for Evidential Reasoning," UCLA CS Technical Report 850021, Proceedings, Cognitive Science Society, UC Irvine, 329-334, August 15-17, 1985.
- [Pearl 00] Pearl, J., Causality: Models, Reasoning, and Inference, Cambridge University Press, 2000
- [Suermondt et al. 90] Henri Jacques Suermondt, Gregory F. Cooper, David Heckerman, "A combination of cutset conditioning with clique-tree propagation in the Pathfinder system", UAI '90.

94

## References

---

- [Pearl 88] Judea Pearl, Probabilistic Reasoning in Intelligent Systems, Morgan Kaufmann: 1998.
- [Lauritzen, Spiegelhalter 1988]
- [Heckerman, Geiger, Chickering 95] D. Heckerman, D. Geiger, D. M. Chickering: "Learning Bayesian Networks: The Combination of Knowledge and Statistical Data", Machine Learning, 20(3), 1995
- [Costa et al 03] V/ Santos Costa, D. Page, M. Qazi, J. Cussens: "CLP(BN): Constraint Logic Programming for Probabilistic Knowledge", UAI03.

95

## References

---

- [Vennekens et al. 04] J.Vennekens, S. Verbaeten and M. Bruynooghe, "Logic programs with annotated disjunctions", ICLP04
- [Riguzzi 07] F. Riguzzi. "A top down interpreter for lpad and cp-logic". In Proceedings of the 10th Congress of the Italian Association for Artificial Intelligence, number 4733 in Lecture Notes in Artificial Intelligence, Springer, 2007.
- [Riguzzi 04] F. Riguzzi. "Learning logic programs with annotated disjunctions", ILP04.

96



## References

---

- [Blooheel, Meert 06] H. Blooheel, W. Meert: “Towards Learning Non-recursive LPADs by Transforming Them into Bayesian Networks”, ILP 2006:
- [Riguzzi 08] F. Riguzzi, “ALLPAD: Approximate learning of logic programs with annotated disjunctions”. Machine Learning, 70(2-3), 2008.
- [Blooheel, Meert 07] W. Meert, J. Struyf, H. Blooheel, “Learning Ground CP-logic Theories by means of Bayesian Network Techniques”, MRDM07