

# Bayesian Networks

Fabrizio Riguzzi

- 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/>

## Summary

---

- Conditional independence
- Definition of Bayesian network
- Inference
- Learning
- Markov networks

# 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

3

# 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 door given that the neighbor called ? (compute  $P(b_2|n_2)$ , belief updating)

4

# 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)

5

# Types of Problems

---

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

6

# 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  $\mathbf{E}$ )

7

## 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  $e, b, a, n$ .

8

# Inference

---

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

$$\begin{aligned} P(\mathbf{q}|\mathbf{e}) &= \frac{P(\mathbf{q}, \mathbf{e})}{P(\mathbf{e})} \\ &= \frac{\sum_{\mathbf{x}, \mathbf{X}=\mathbf{U}\setminus\mathbf{Q}\setminus\mathbf{E}} P(\mathbf{x}, \mathbf{q}, \mathbf{e})}{\sum_{\mathbf{y}, \mathbf{Y}=\mathbf{U}\setminus\mathbf{E}} P(\mathbf{y}, \mathbf{e})} \end{aligned}$$

9

# Computational Cost

---

- If we have  $n$  binary variables ( $|\mathbf{U}|=n$ ), knowing the jpd requires storing  $O(2^n)$  different values.
- Even if we 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

10

# Conditional Independence

---

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

$$\forall \mathbf{x}, \mathbf{y}, \mathbf{z}: P(\mathbf{y}, \mathbf{z}) > 0 \rightarrow P(\mathbf{x}|\mathbf{y}, \mathbf{z}) = P(\mathbf{x}|\mathbf{z})$$

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

$$\forall \mathbf{x}, \mathbf{y}: P(\mathbf{y}) > 0 \rightarrow P(\mathbf{x}|\mathbf{y}) = P(\mathbf{x})$$

- We write  $I\langle \mathbf{X}, \{\}, \mathbf{Y} \rangle$

11

---

# 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:

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

12

# 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  $\Pi_i$   
$$P(x_i | x_{i-1}, \dots, x_1) = P(x_i | \pi_i) \quad \text{whenever } P(x_{i-1}, \dots, x_1) > 0$$
- where  $\pi_i$  is a set of values for  $\Pi_i$
- $\Pi_i$  parents of  $X_i$

13

# Conditional Independence

---

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

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

14

# 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  $O(2^n)$

15

# Graphical Representation

---

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

16



# Example Network

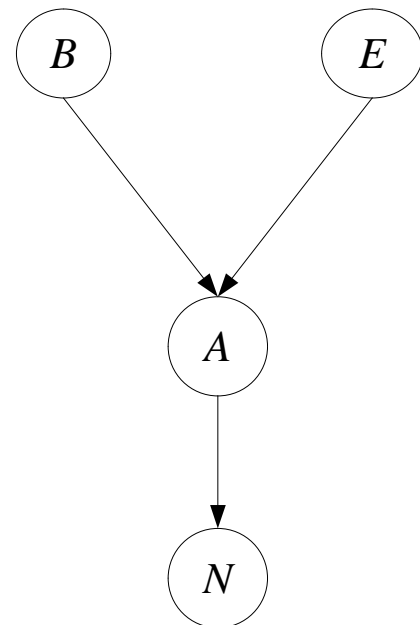
- Variable order: E,B,A,N
- Independencies

$$P(e)$$

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

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

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



17

## Conditional Probability Tables

- Earthquake E,  $e_1$ =no,  $e_2$ =moderate,  $e_3$ =severe
- Burglary B, :  $b_1$ =no,  $b_2$ =yes through door,  $b_3$ =yes through window
- Alarm A,  $a_1$ =no,  $a_2$ =yes
- Neighbor call N,  $n_1$ =no,  $n_2$ =yes

P(B)	
B=no	0,7
B=door	0,1
B=windows	0,2

P(E)	
E=no	0,6
E=moderate	0,2
E=severe	0,2

P(A EB)	no,no	no,do	no,wi	mo,no	mo,do	mo,wi	se,no	se,do	se,wi
no	0,99	0,1	0,2	0,8	0,08	0,1	0,7	0,05	0,07
yes	0,01	0,9	0,8	0,2	0,92	0,9	0,3	0,95	0,93

P(N A)	A=no	A=yes
N=no	0,9	0,05
N=yes	0,1	0,95

18

# 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's)
$$\{\theta_{x_i|\pi_i} \in R \mid i = 1, \dots, n, x_i \in X_i, \pi_i \in \Pi_i\}$$
  - where  $\Pi_i$  is the set of parents of  $X_i$

19

# Bayesian Network

---

- BNs are also called belief networks or directed acyclic graphical models

20

# Bayesian Network

---

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

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

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

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

21

---

## 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  $\pi_i$

22

# Building a Bayesian Network

---

- Usually the expert considers 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]

23

## 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.

24

## Pathfinder system [Suermondt et al. 90]

---

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

25

## How to Tell Independence

---

- There is a relatively simple algorithm for determining whether two variables in a Bayesian network are conditionally independent: **d-separation**.
- Definition:  $X$  and  $Z$  are **d-separated** by a set of evidence variables  $E$  iff every undirected path from  $X$  to  $Z$  is “blocked”, where a path is “blocked” iff one or more of the following conditions is true: ...

26

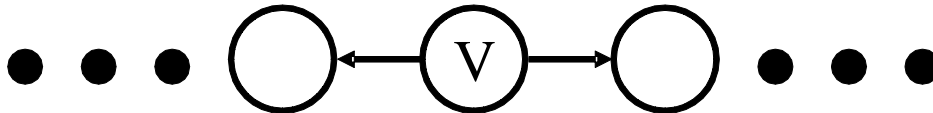
# Blocked Path

---

There exists a variable  $V$  on the path such that

it **is** in the evidence set  $E$

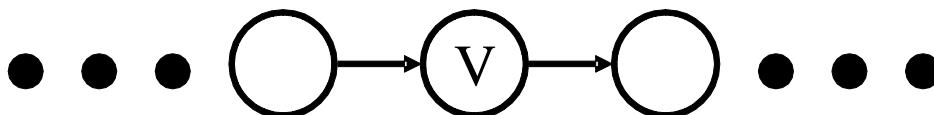
the arcs putting  $V$  in the path are “tail-to-tail”



Or, there exists a variable  $V$  on the path such that

it **is** in the evidence set  $E$

the arcs putting  $V$  in the path are “tail-to-head”



27

# Blocked Path

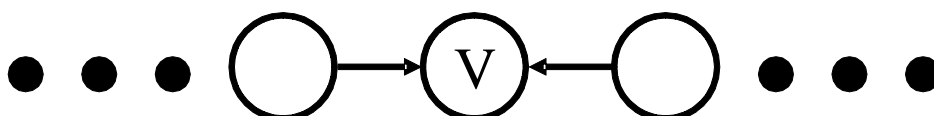
---

- ... Or, there exists a variable  $V$  on the path such that

it **is NOT** in the evidence set  $E$

**neither are any of its descendants**

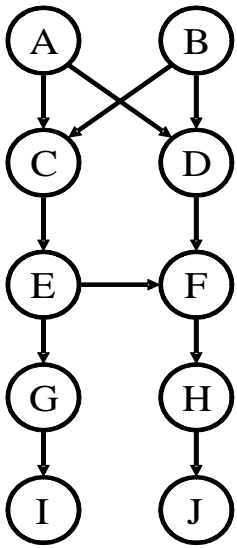
the arcs putting  $V$  on the path are “head-to-head”



28

## Example

---

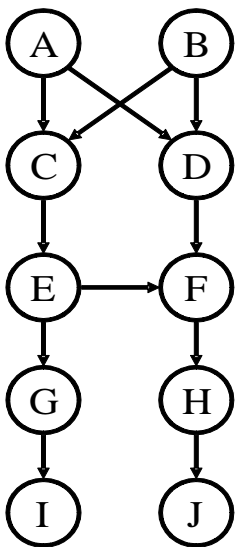


- $I\langle C, \{\}, D\rangle?$
- $I\langle C, \{A\}, D\rangle?$
- $I\langle C, \{A, B\}, D\rangle?$
- $I\langle C, \{A, B, J\}, D\rangle?$
- $I\langle C, \{A, B, E, J\}, D\rangle?$

29

## Example

---



- $I\langle C, \{\}, D\rangle?$  No
- $I\langle C, \{A\}, D\rangle?$  No
- $I\langle C, \{A, B\}, D\rangle?$  Yes
- $I\langle C, \{A, B, J\}, D\rangle?$  No
- $I\langle C, \{A, B, E, J\}, D\rangle?$  Yes

30

# Inference with Bayesian Networks

---

- With a Bayesian Network we save space, do we also save time?

- Do we have to use the formula

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

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

31

# 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 simulation
    - Loopy belief propagation
    - Variational methods,

32



# Complexity of Inference

---

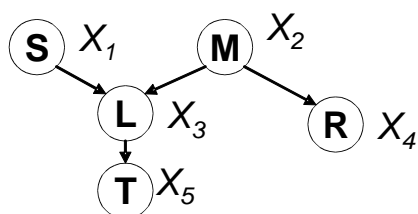
- Exact inference with BN is #P-complete
- #P-complete: a special case of NP-complete problems
  - The answer to a #P-complete problem is the number of solutions to a NP-complete problem

33

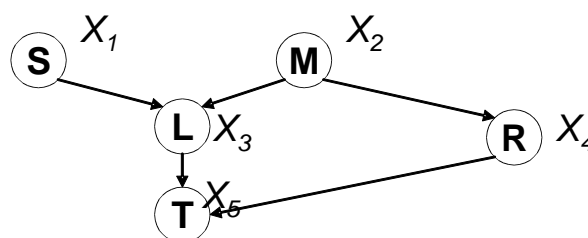
## Polytrees

---

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

34

# Belief Propagation [Pearl 88]

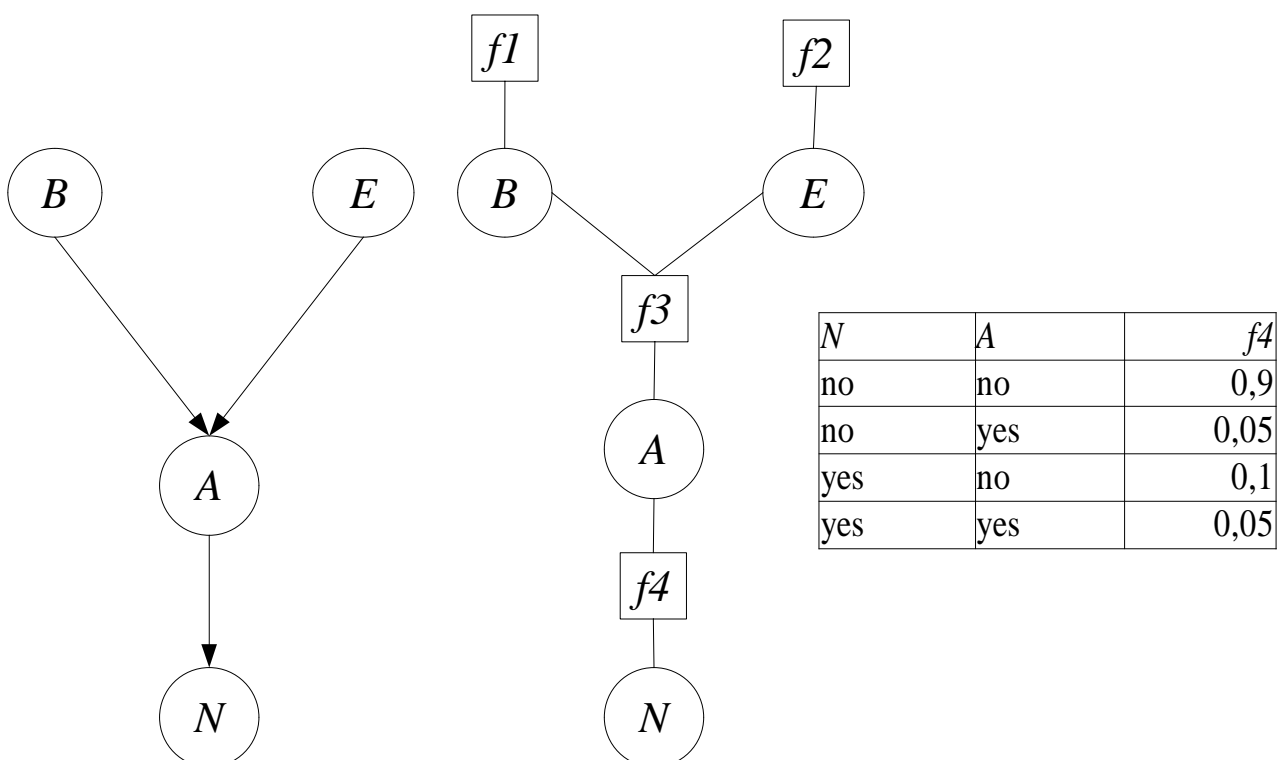
---

- Best presented over Factor Graphs
- A **Factor Graph** is a bipartite graph  $(V,F,E)$  where vertices  $V$  index the variables, the vertices  $F$  index the families (factors), and edges  $E$  are connected between  $V$  and  $F$
- A factor, given the values of the variables involved in the factor, returns a non-negative number.
- A family in a BN can be seen as a factor

35

## Example Network

---



36

# Messages

---

- The message from a variable node  $X$  to a neighbor factor node  $f$  is

$$\mu_{X \rightarrow f}(x) = \prod_{h \in \text{nb}(X) \setminus X} \mu_{h \rightarrow X}(x)$$

- where  $\text{nb}(X)$  is the set of neighbor of  $X$ , the set of factors  $X$  appears in
- The message from a factor to a variable is

$$\mu_{f \rightarrow X}(x) = \sum_{\neg\{X\}} (f(\mathbf{x}) \prod_{Y \in \text{nb}(f) \setminus X} \mu_{Y \rightarrow f}(y))$$

- Where  $\text{nb}(f)$  is the set of arguments of  $f$  and the sum is over all of these except  $X$

37

# Belief

---

- The unnormalized belief of each variable  $X_i$  in iteration  $k$  can be computed from the equation

$$b_i(x_i) = \prod_{f \in \text{nb}(X_i)} \mu_{f \rightarrow X_i}(x_i)$$

- For example, if  $X_1$  has 3 values  $x_{11}$ ,  $x_{12}$ ,  $x_{13}$ , their probabilities are
- $B = b_1(x_{11}) + b_1(x_{12}) + b_1(x_{13})$
- $P(x_{11}) = b_1(x_{11})/B$      $P(x_{12}) = b_1(x_{12})/B$      $P(x_{13}) = b_1(x_{13})/B$

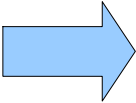
38

# Incorporation of Evidence

---

- For each factor  $f$ , for each combination of values of the arguments that is incompatible with the evidence,  $f(\mathbf{x})$  is set to 0
- Example: evidence  $N=\text{yes}$ , factor  $f_4$  becomes

$N$	$A$	$f_4$
no	no	0,9
no	yes	0,05
yes	no	0,1
yes	yes	0,05



$N$	$A$	$f_4$
no	no	0
no	yes	0
yes	no	0,1
yes	yes	0,05

39

# Algorithm

---

- Initialize all messages to 1 or randomly
- Loop
  - Select an arc
  - Compute the value of the message on the arc
- Until the messages do not change anymore
- If the network is a polytree, this algorithm converges
- Various strategies for selecting the arc to update

40

# Message schedules

---

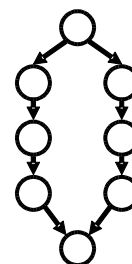
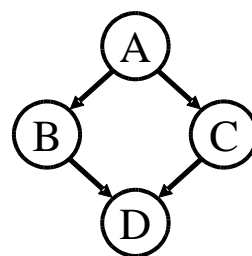
- The order in which messages are updated
- Asynchronous schedules: messages are updated sequentially, one arc at a time
- Synchronous schedules: all messages are updated in parallel.
- Flooding (asynchronous): messages are passed from each variable to each corresponding factor and back at each step
- The most widely used and generally best-performing method

41

# General Networks

---

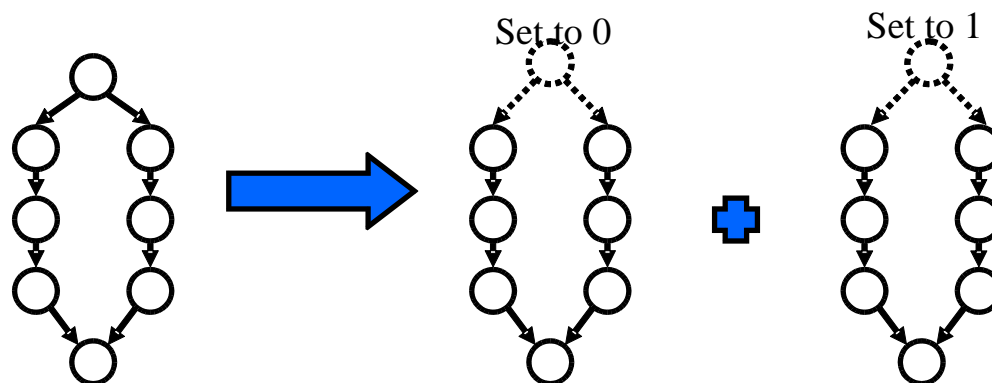
- Networks that have a cycle in their undirected version
- Three possibilities
  - Conditioning
  - Clustering
  - Approximations



42

# Conditioning

---

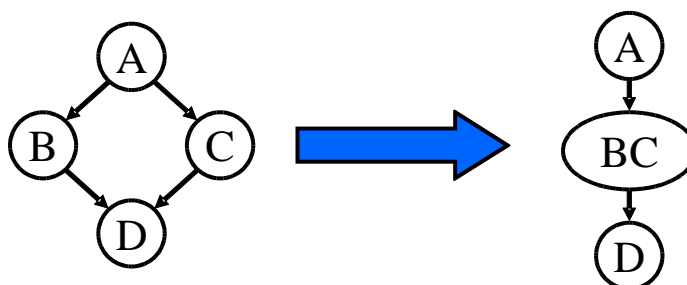


43

# Clustering

---

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



- Problem: how to find a good clustering?

44

# 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

45

# Junction Tree

---

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

46

# Approximate Methods

---

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

47

# Stochastic Simulation

---

- 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

48



## 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

49

## 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  $D$ ?

50

# 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

$$P(D) = \prod_{i=1}^N P(\mathbf{u}^i)$$

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

$$\log P(D) = \sum_{i=1}^N \log P(\mathbf{u}^i)$$

51

# 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  $\mathbf{u}^j$ , the values of all the variables are observed
  - incomplete data

52

# 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

53

## 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

$$C_{\pi_i} = \langle N_{x_i^1, \pi_i}, \dots, N_{x_i^{v_i}, \pi_i} \rangle$$

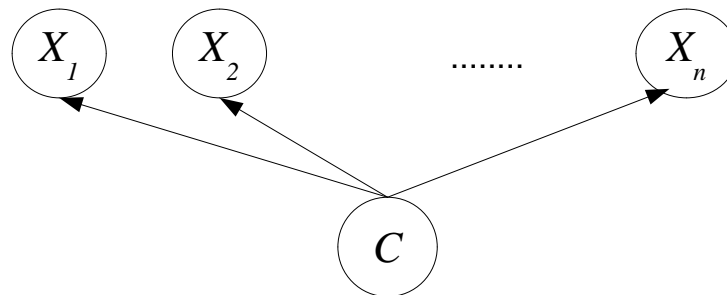
- where  $v_i$  is the number of values of  $X_i$

54

# Naive Bayes Special Case

---

- We want to perform classification
- One variable  $C$  represents the class
- The variables  $\mathbf{X}$  represent the attributes
- Model:

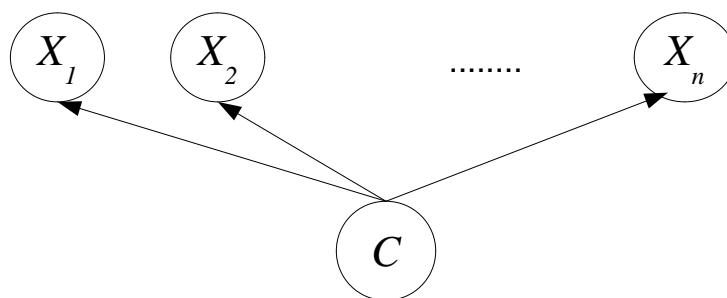


- $X_i$  independent from  $X_j$  given  $C$

55

# Naive Bayes Special Case

---



- Conditional probability tables (case of Boolean variables):

	C=true	C=false
$X_i$ =true	$P(X_i=\text{true} C=\text{true})$	$P(X_i=\text{true} C=\text{false})$
$X_i$ =false	$P(X_i=\text{false} C=\text{true})$	$P(X_i=\text{false} C=\text{false})$

56

# Example

No	Outlook	Temp	Humid	Windy	Class
D1	sunny	mild	normal	T	P
D2	sunny	hot	high	T	N
D3	sunny	hot	high	F	N
D4	sunny	mild	high	F	N
D5	sunny	cool	normal	F	P
D6	overcast	mild	high	T	P
D7	overcast	hot	high	F	P
D8	overcast	cool	normal	T	P
D9	overcast	hot	normal	F	P
D10	rain	mild	high	T	N
D11	rain	cool	normal	T	N
D12	rain	mild	normal	F	P
D13	rain	cool	normal	F	P
D14	rain	mild	high	F	P

	C=P	C=N
Humid=normal	6/9=0.66666	1/5=0.2
Humid=high	3/9=0.33333	4/5=0.8

57

# Queries

- Computing the probability of a class given values for the attributes:  $P(c|x_1, \dots, x_n)$

$$P(c|x_1, \dots, x_n) = \frac{P(c, x_1, \dots, x_n)}{P(x_1 \dots x_n)} = \frac{P(x_1, \dots, x_n|c) P(c)}{P(x_1 \dots x_n)}$$

- Since the attributes are independent given the class

$$P(c|x_1, \dots, x_n) = \frac{P(x_1|c) \dots P(x_n|c) P(c)}{P(x_1 \dots x_n)}$$

58

# Example

---

- We want to classify  $\langle \text{Outlook}=\text{sunny}, \text{Temp}=\text{cool}, \text{Humid}=\text{high}, \text{Windy}=\text{T} \rangle$
- We have to compute  $P(\text{Class}=\text{P} | \text{Outlook}=\text{sunny}, \text{Temp}=\text{cool}, \text{Humid}=\text{high}, \text{Windy}=\text{T})$
- We compute only the parameters we need

$$P(\text{Class}=\text{P})=9/14=0.64$$

$$P(\text{Class}=\text{N})=5/14=0.36$$

$$P(\text{Outlook}=\text{sunny} | \text{Class}=\text{P})=2/9=0.222$$

$$P(\text{Outlook}=\text{sunny} | \text{Class}=\text{N})=3/5=0.6$$

$$P(\text{Temp}=\text{cool} | \text{Class}=\text{P})=3/9=0.333$$

$$P(\text{Temp}=\text{cool} | \text{Class}=\text{N})=1/5=0.2$$

$$P(\text{Humid}=\text{high} | \text{Class}=\text{P})=3/9=0.333$$

$$P(\text{Humid}=\text{high} | \text{Class}=\text{N})=4/5=0.8$$

$$P(\text{Windy}=\text{T} | \text{Class}=\text{P})=3/9=0.33$$

$$P(\text{Windy}=\text{T} | \text{Class}=\text{N})=3/5=0.6$$

59

# Example

---

$$P(\text{Class}=\text{P}, \text{Outlook}=\text{sunny}, \text{Temp}=\text{cool}, \text{Humid}=\text{high}, \text{Windy}=\text{T}) = 0.0053$$

$$P(\text{Class}=\text{N}, \text{Outlook}=\text{sunny}, \text{Temp}=\text{cool}, \text{Humid}=\text{high}, \text{Windy}=\text{T}) = 0.0206$$

- We can compute  $P(\text{Outlook}=\text{sunny}, \text{Temp}=\text{cool}, \text{Humid}=\text{high}, \text{Windy}=\text{T})$  by marginalization:

$$P(\text{Outlook}=\text{sunny}, \text{Temp}=\text{cool}, \text{Humid}=\text{high}, \text{Windy}=\text{T}) =$$

$$P(\text{Class}=\text{P}, \text{Outlook}=\text{sunny}, \text{Temp}=\text{cool}, \text{Humid}=\text{high}, \text{Windy}=\text{T}) +$$

$$P(\text{Class}=\text{N}, \text{Outlook}=\text{sunny}, \text{Temp}=\text{cool}, \text{Humid}=\text{high}, \text{Windy}=\text{T}) =$$

$$0.0053 + 0.0206 = 0.0259$$

- So

$$P(\text{Class}=\text{P} | \text{Outlook}=\text{sunny}, \text{Temp}=\text{cool}, \text{Humid}=\text{high}, \text{Windy}=\text{T}) = 0.0053 / 0.0259 = 0.205$$

$$P(\text{Class}=\text{N} | \text{Outlook}=\text{sunny}, \text{Temp}=\text{cool}, \text{Humid}=\text{high}, \text{Windy}=\text{T}) = 0.0206 / 0.0259 = 0.795$$

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

---

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

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

- where

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

- 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  $X_i$

$$C_i = \langle C_{\pi_i^1}, \dots, C_{\pi_i^{r_i}} \rangle$$

$$C'_i = \langle C'_{\pi_i^1}, \dots, C'_{\pi_i^{r_i}} \rangle$$

64



# Structure Score

---

$$BD_i(G) = \sum_{\pi_i} \log \frac{\Gamma(N_{\pi_i})}{\Gamma(N_{\pi_i} + N'_{\pi_i})} + \sum_{x_i} \log \frac{\Gamma(N_{x_i, \pi_i} + N'_{x_i, \pi_i})}{\Gamma(N_{x_i, \pi_i})}$$

- Where  $\Gamma$  is the Gamma function, an extension of the factorial function with its argument shifted down by 1, to real and complex numbers. That is, if  $n$  is a positive integer:

$$\Gamma(n) = (n-1)!$$

- otherwise

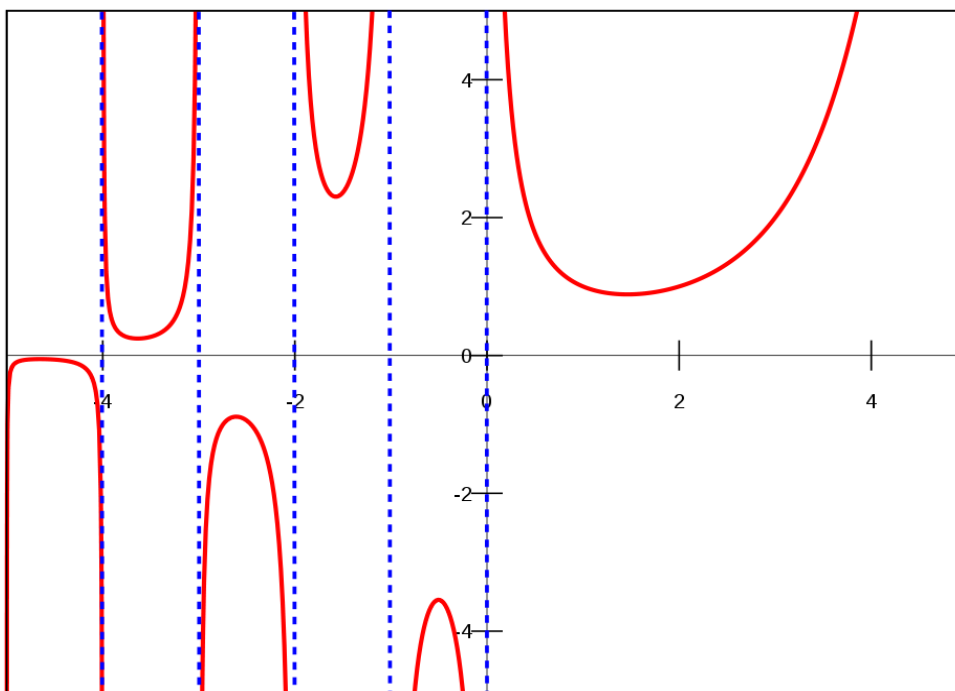
$$\Gamma(z) = \int_0^{\infty} t^{z-1} e^{-t} dt$$

65

# Gamma Function

---

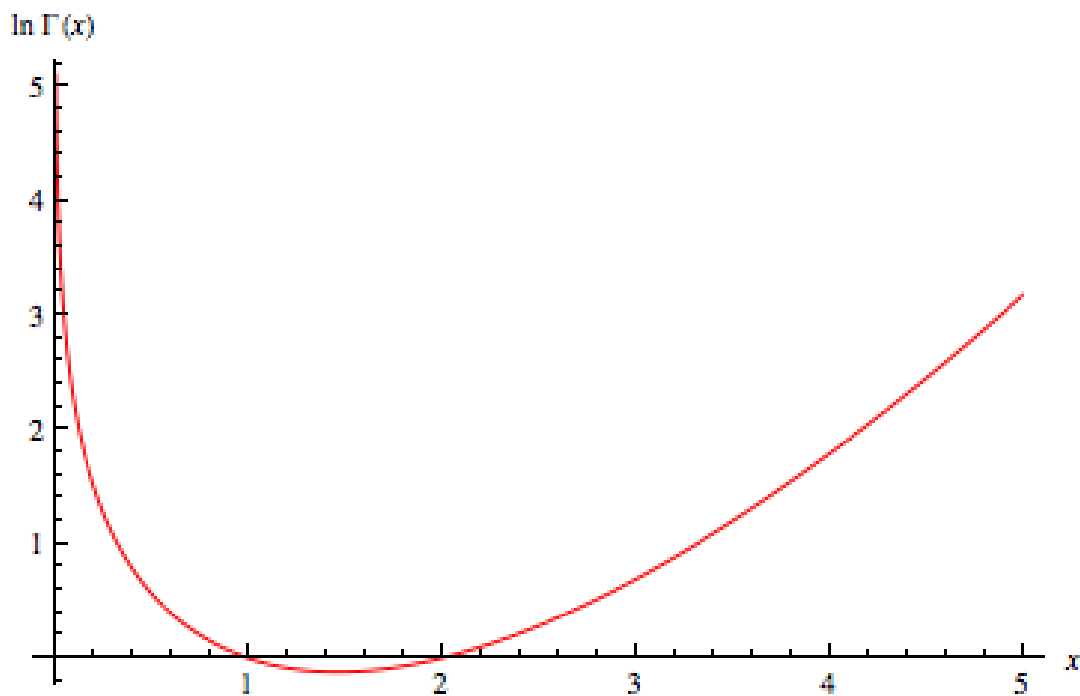
Gamma function



66

# Log Gamma Function

---



67

## Structure Score

---

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

68

# 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

69

## EM Algorithm

---

- Initialize the parameters at random  $\Theta$
- Repeat
  - Expectation step:
    - compute the probability  $P(y|e)$  of each value  $y$  of the missing attributes using  $(G, \Theta)$  and inference
  - Compute  $\Theta$  by maximum likelihood on  $D'$ 
    - Relative frequency for each family
    - If a variable  $Y$  is unobserved in an example  $e$  that matches  $x_i, \pi_i$ , instead of adding 1 to  $N_{x_i, \pi_i}$  we add  $P(y|e)$

70

# Structure Learning from Incomplete Data

---

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

71

## 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

72

# Applications of BN

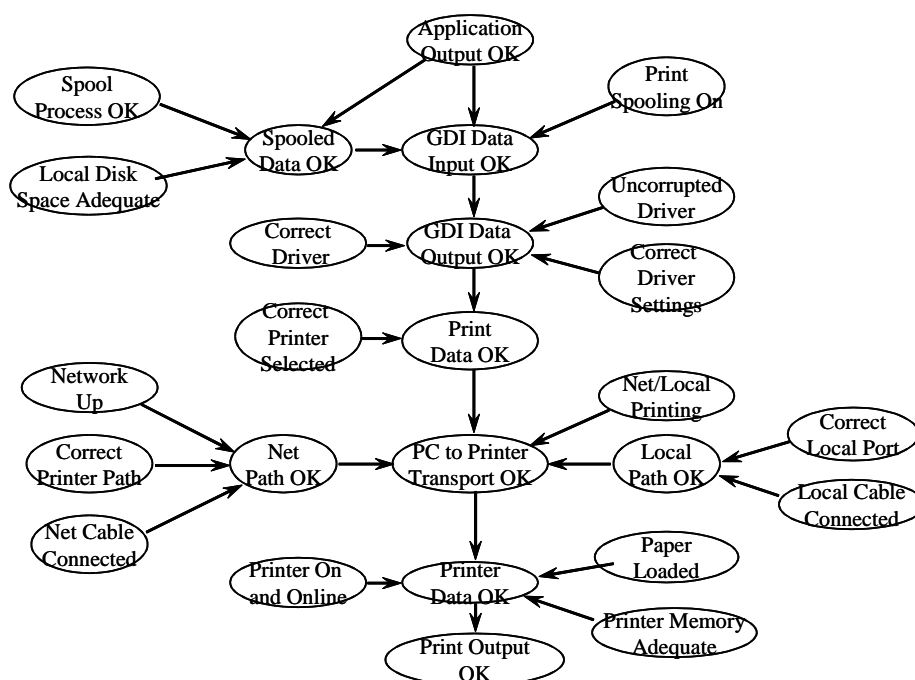
---

- 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)

73

## Printer Troubleshooting (Windows 95)

---



74

# 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

75

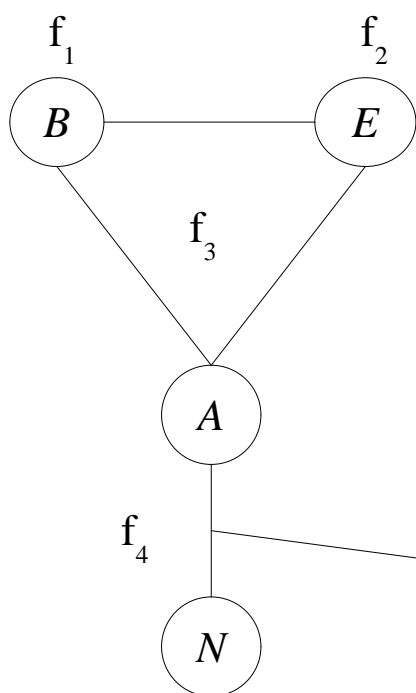
## Markov Networks (MN)

---

- Approach alternative to BN
- Also called Random Fields, undirected graphical models
- Undirected graph
- Conditional independence represented by graph separation
- Probability distribution as the product of a set of **potentials** or **factors** (functions of a subset of variables) divided by a normalization constant
- Potentials over cliques

76

# Example



- Four potentials  $f_1$ ,  $f_2$ ,  $f_3$  and  $f_4$

$N$	$A$	$f_4$
no	no	0,9
no	yes	0,05
yes	no	0,1
yes	yes	0,05

77

# Markov Networks

- Probability

$$P(\mathbf{u}) = \frac{\prod_c f_c(\mathbf{x}_c)}{Z}$$

$$Z = \sum_u \prod_c f_c(\mathbf{x}_c)$$

- $Z$  is called **partition function**, ensures that the probabilities sum to 1

78

# Loglinear Models

---

- If all the potentials are  $>0$ , they can be represented as exponential functions, i.e.,  $f_4$  can be represented as  $f_4 = \exp(w_4 F_4)$
- where  $F_4$  is any real function of  $f_4$  arguments and  $w_4$  is a real weight. Then

$$P(\mathbf{u}) = \frac{\exp \sum_c w_c F_c(\mathbf{x}_c)}{Z}$$

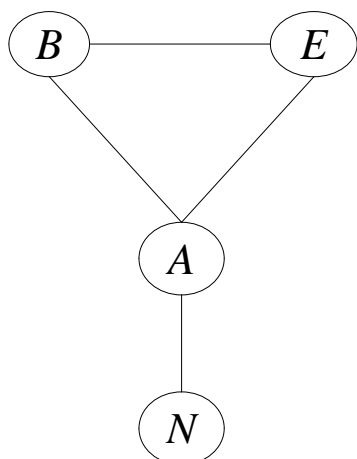
$$Z = \sum_u \exp \sum_c w_c F_c(\mathbf{x}_c)$$

79

## How to tell independence

---

- Definition:  $\mathbf{X}$  and  $\mathbf{Y}$  are **independent** given a set of variables  $\mathbf{Z}$  ( $I\langle \mathbf{X}, \mathbf{Z}, \mathbf{Y} \rangle$ ) iff every path from  $\mathbf{X}$  to  $\mathbf{Y}$  passes through a variable of  $\mathbf{Z}$



- $I\langle B, \{\}, N \rangle$ ?
- $I\langle B, A, N \rangle$ ?
- $I\langle B, E, N \rangle$ ?
- $I\langle \{B, E\}, A, N \rangle$ ?

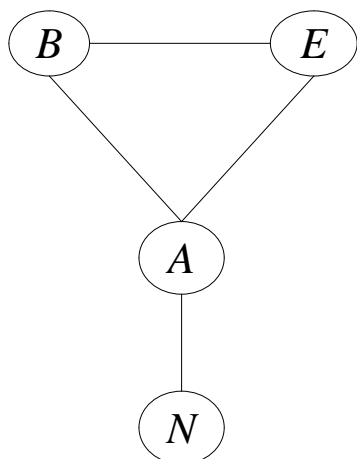
80



# How to tell independence

---

- Definition:  $\mathbf{X}$  and  $\mathbf{Y}$  are **independent** given a set of variables  $\mathbf{Z}$  ( $I\langle\mathbf{X},\mathbf{Z},\mathbf{Y}\rangle$ ) iff every path from  $\mathbf{X}$  to  $\mathbf{Y}$  passes through a variable of  $\mathbf{Z}$



- $I\langle B, \{\}, N \rangle$  No
- $I\langle B, A, N \rangle$  Yes
- $I\langle B, E, N \rangle$  No
- $I\langle \{B, E\}, A, N \rangle$  Yes

81

# Markov Network

---

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

82

# 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

83

# Resources

---

- Daphne Koller, Nir Friedman, Probabilistic graphical models: principles and techniques, MIT Press: 2009, ISBN 978-0-262-01319-2
- 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/>

84

## 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.

85

## 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

86