

Bayesian AI Tutorial

Ann E. Nicholson and Kevin B. Korb

Faculty of Information Technology
Monash University
Clayton, VIC 3168
AUSTRALIA

{annn,korb}@csse.monash.edu.au

[HTTP://WWW.CSSE.MONASH.EDU.AU/BAI](http://www.csse.monash.edu.au/bai)

Text: *Bayesian Artificial Intelligence*, Kevin B. Korb
and Ann E. Nicholson, Chapman & Hall/CRC, 2004.

Overview

1. Bayesian AI
2. Introduction to Bayesian networks
3. Extensions to Bayesian networks
 - (a) Decision networks
 - (b) Dynamic Bayesian networks
4. Learning Bayesian networks
5. Knowledge Engineering with Bayesian networks
6. Monash BAI Group: BN Applications

Additional material

- References: (p. XXX)
- *Bayesian Artificial Intelligence* Table of Contents
- *Bayesian Artificial Intelligence* Appendix B:
Software Packages

Introduction to Bayesian AI

- Reasoning under uncertainty
- Probabilities
- Bayesian philosophy
 - Bayes' Theorem
 - Conditionalization
 - Motivation
- Bayesian decision theory
- How to be an effective Bayesian
- Probabilistic causality
 - Humean causality
 - Prob causality
 - Are Bayesian networks Bayesian?
- Towards a Bayesian AI

Reasoning under uncertainty

Uncertainty: The quality or state of being not clearly known.

This encompasses most of what we understand about the world — and most of what we would like our AI systems to understand.

Distinguishes *deductive* knowledge (e.g., mathematics) from *inductive* belief (e.g., science).

Sources of uncertainty

- Ignorance
(which side of this coin is up?)
- Complexity
(meteorology)
- Physical randomness
(which side of this coin will land up?)
- Vagueness
(which tribe am I closest to genetically? Picts? Angles? Saxons? Celts?)

Probability Calculus

Classic approach to reasoning under uncertainty.
(origin: Blaise Pascal and Fermat).

Kolmogorov's Axioms:

1. $P(U) = 1$

2. $\forall X \subseteq U \ P(X) \geq 0$

3. $\forall X, Y \subseteq U$

if $X \cap Y = \emptyset$

then $P(X \vee Y) = P(X) + P(Y)$

Conditional Probability $P(X|Y) = \frac{P(X \wedge Y)}{P(Y)}$

Independence $X \perp\!\!\!\perp Y$ iff $P(X|Y) = P(X)$

Rev. Thomas Bayes
(1702-1761)



Bayesian AI Tutorial

Bayes' Theorem; Conditionalization

— Due to Reverend Thomas Bayes (1764)

$$P(h|e) = \frac{P(e|h)P(h)}{P(e)}$$

Conditionalization: $P'(h) = P(h|e)$

Or, read Bayes' theorem as:

$$\text{Posterior} = \frac{\text{Likelihood} \times \text{Prior}}{\text{Prob of evidence}}$$

Assumptions:

1. Joint priors over $\{h_i\}$ and e exist.
2. Total evidence: e , and only e , is learned.

Motivation: Breast Cancer

Let $P(h) = 0.01$ (one in 100 women tested have it)

$P(e|h) = 0.8$ and $P(e|\neg h) = 0.1$

(true and false positive rates). What is $P(h|e)$?

Motivation: Breast Cancer

Let $P(h) = 0.01$ (one in 100 women tested have it)

$P(e|h) = 0.8$ and $P(e|\neg h) = 0.1$

(true and false positive rates). What is $P(h|e)$?

Bayes' Theorem yields:

$$\begin{aligned} P(h|e) &= \frac{P(e|h)P(h)}{P(e)} \\ &= \frac{P(e|h)P(h)}{P(e|h)P(h) + P(e|\neg h)P(\neg h)} \\ &= \frac{0.8 \times 0.01}{0.8 \times 0.01 + 0.1 \times 0.99} \\ &= \frac{0.008}{0.008 + 0.099} \\ &= \frac{0.008}{0.107} \\ &\approx 0.075 \end{aligned}$$

Motivation

Huge variety of cases where

- Uncertainty dominates considerations
- Getting it right is crucial

Examples and consequences:

- Medicine: death, injury, disease
- Law: false imprisonment, wrongful execution
- Space shuttle: explosion
- Hiring: wasted time and money

Humean Causality

As we shall see, causality and Bayesian networks are intimately related concepts.

David Hume (1737) analyzed *A causes B* as:

- Whenever *A* occurs, *B* occurs
- *A* and *B* are contiguous
- *A* precedes *B*

This was immediately challenged by Thomas Reid: let *A* be night and *B* day; the conditions are satisfied, but neither causes the other.

Leading to: CounterExample → new conditions → CE ...

Through the next two centuries the “conditions” (sufficiency) account of causality has built up complexity without explanation

Probabilistic Causality

Salmon (1980): What is this sufficiency nonsense?

Either of determinism and indeterminism are *possible* – i.e., it is a contingent fact of the world whether it is deterministic or not.

1. A philosophical analysis of causality should not *presume* determinism.
2. Besides, there is evidence that indeterminism is correct.
3. A probabilistic analysis of causality does predetermine the determinism question, whereas the sufficiency analysis does.

Probabilistic causality

- started by H Reichenbach, IJ Good, P Suppes, W Salmon
- has turned into the theory of Bayesian networks

Bayesian AI

A Bayesian conception of an AI is:

An autonomous agent which

- Has a utility structure (preferences)
- Can learn about its world and the relation between its actions and future states (probabilities)
- Maximizes its expected utility

The techniques used in learning about the world are (primarily) statistical... Hence

Bayesian data mining

Bayesian Networks: Overview

- Introduction to BNs
 - Nodes, structure and probabilities
 - Reasoning with BNs
 - Understanding BNs
- Inference
 - Exact algorithms
 - Approximate algorithms
 - Causal modelling
- Extensions to BNs
 - Decision networks
 - Dynamic Bayesian networks (DBNs)

Bayesian Networks

- Data Structure which represents the dependence between variables.
- Gives concise specification of the joint probability distribution.
- A Bayesian Network is a graph in which the following holds:
 1. A set of random variables makes up the nodes in the network.
 2. A set of directed links or arrows connects pairs of nodes.
 3. Each node has a conditional probability table that *quantifies* the effects the parents have on the node.
 4. Directed, acyclic graph (DAG), i.e. no directed cycles.

Example: Lung Cancer Diagnosis

A patient has been suffering from shortness of breath (called dyspnoea) and visits the doctor, worried that he has lung cancer. The doctor knows that other diseases, such as tuberculosis and bronchitis are possible causes, as well as lung cancer. She also knows that other relevant information includes whether or not the patient is a smoker (increasing the chances of cancer and bronchitis) and what sort of air pollution he has been exposed to. A positive XRay would indicate either TB or lung cancer.

Nodes and Values

Q: What are the nodes to represent and what values can they take?

Nodes can be discrete or continuous; will focus on discrete for now.

- Boolean nodes: represent propositions, taking binary values true (T) and false (F).

Example: *Cancer* node represents proposition “the patient has cancer”.

- Ordered values..

Example: *Pollution* node with values $\{low, medium, high\}$.

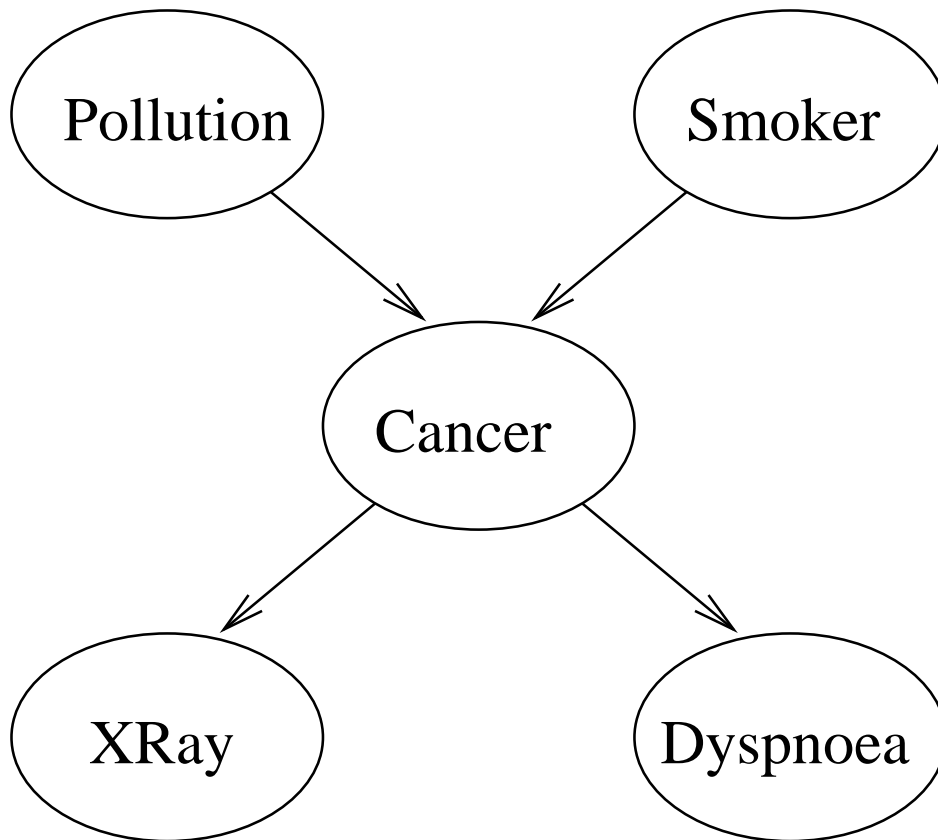
- Integral values.

Example: *Age* node with possible values from 1 to 120.

Lung cancer example: nodes and values

Node name	Type	Values
<i>Pollution</i>	Binary	{ <i>low, high</i> }
<i>Smoker</i>	Boolean	{ <i>T, F</i> }
<i>Cancer</i>	Boolean	{ <i>T, F</i> }
<i>Dyspnoea</i>	Boolean	{ <i>T, F</i> }
<i>XRay</i>	Binary	{ <i>pos, neg</i> }

Lung cancer example: network structure



Note: No explicit representation of other causes of cancer, or other causes of symptoms.

Structure terminology and layout

- Family metaphor:

Parent \Rightarrow *Child*

Ancestor $\Rightarrow \dots \Rightarrow$ *Descendant*

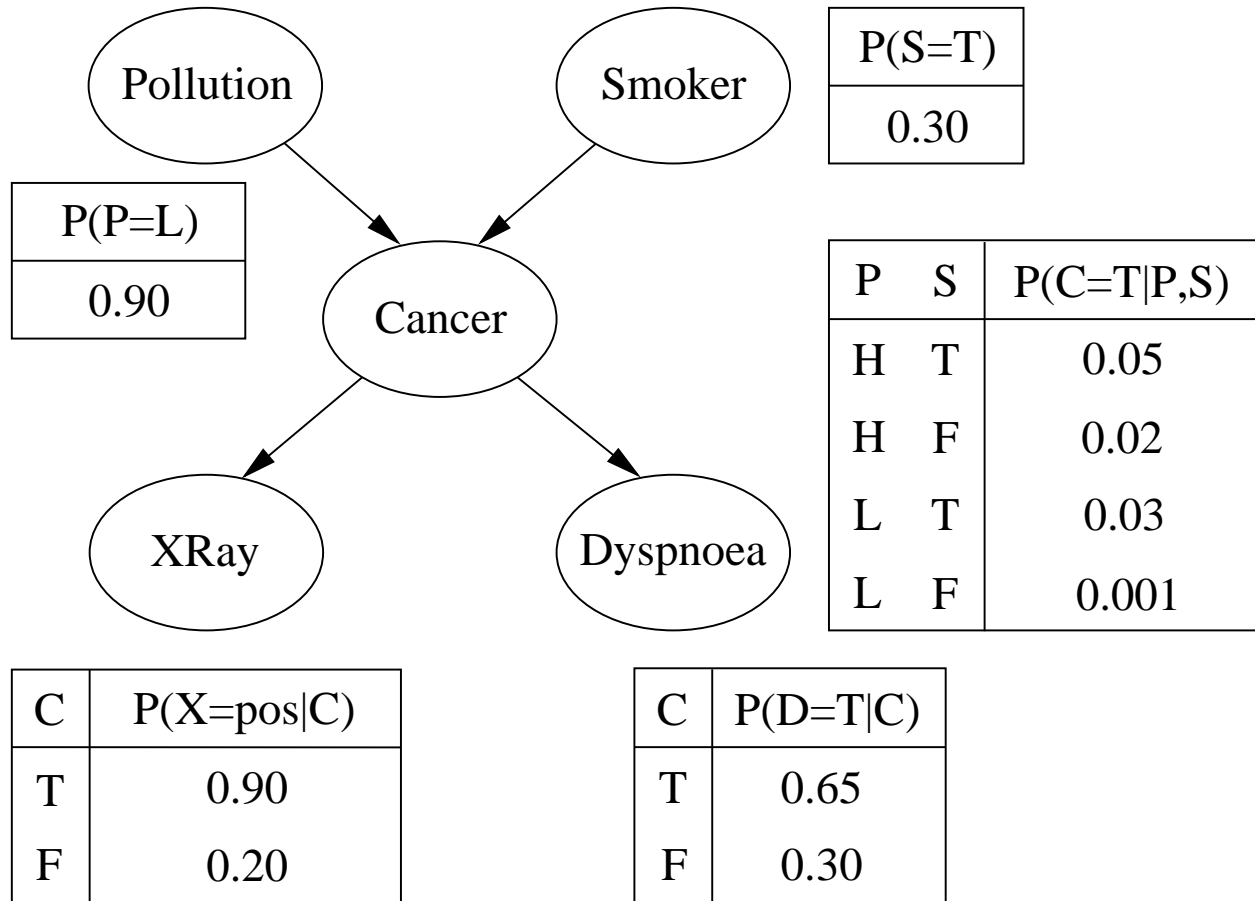
- Markov Blanket = parents + children + children's parents
- Tree analogy:
 - **root** node: no parents
 - **leaf** node: no children
 - **intermediate** node: non-leaf, non-root
- Layout convention: root nodes at top, leaf nodes at bottom, arcs point down the page.

Conditional Probability Tables

Once specified topology, need to specify **conditional probability table** (CPT) for each node.

- Each row contains the conditional probability of each node value for a each possible combination of values of its parent nodes.
- Each row must sum to 1.
- A table for a Boolean var with n Boolean parents contain 2^{n+1} probs.
- A node with no parents has one row (the prior probabilities)

Lung cancer example: CPTs



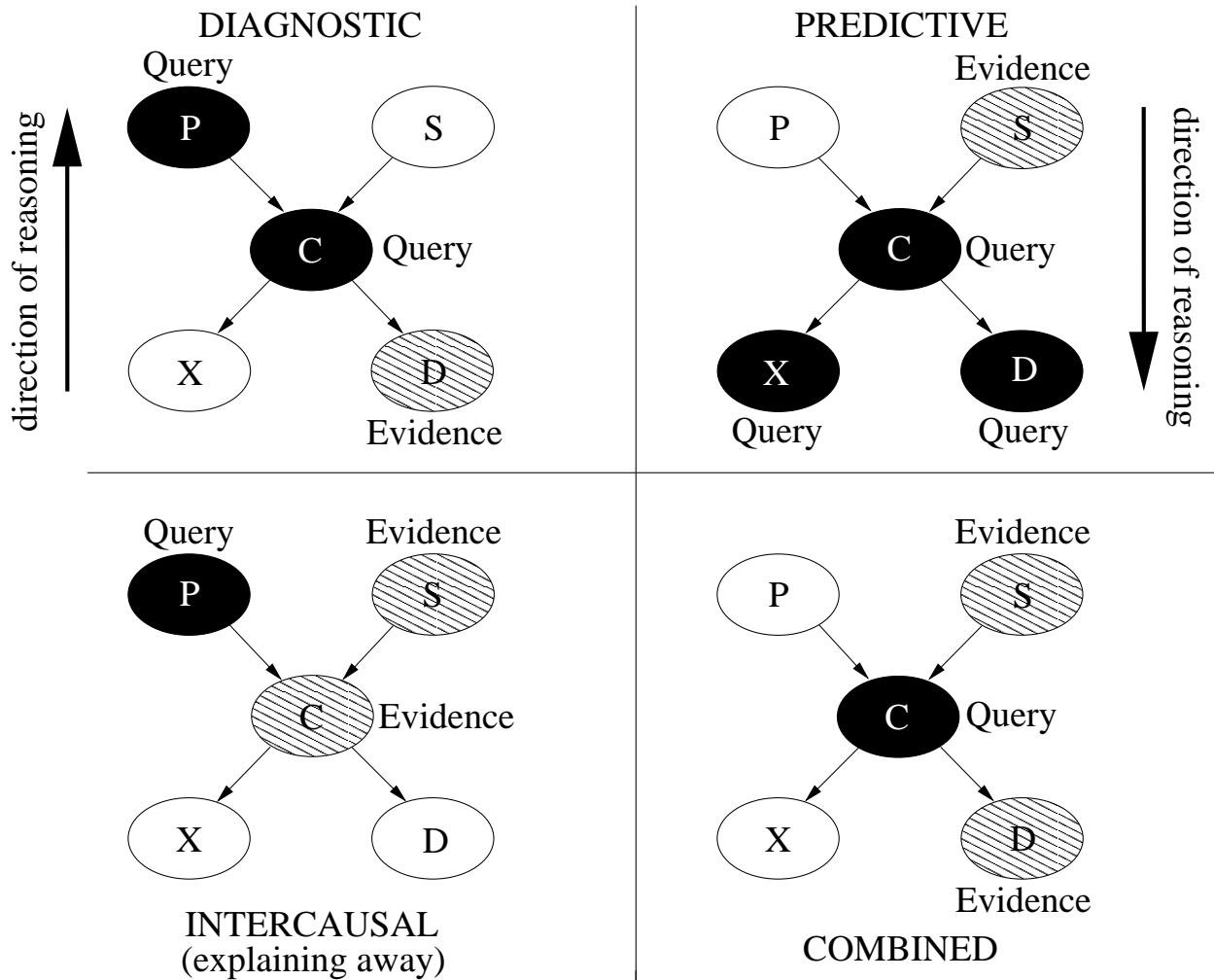
The Markov Property

- Modelling with BNs requires the assumption of the **Markov Property**:
there are no direct dependencies in the system being modeled which are not already explicitly shown via arcs.
- Example: there is no way for smoking to influence dyspnoea except by way of causing cancer.
- BNs which have the Markov property are called Independence-Maps (I-Maps).
- Note: existence of arc does not have to correspond to real dependency in the system being modelled - can be nullified in the CPT.

Reasoning with Bayesian Networks

- Basic task for any probabilistic inference system:
Compute the posterior probability distribution for a set of **query variables**, given new information about some **evidence variables**.
- Also called *conditioning* or *belief updating* or *inference*.

Types of Reasoning



Types of Evidence

- Specific evidence: a definite finding that a node X has a particular value, x .

Example: $Smoker=T$

- Negative evidence: a finding that node Y is *not* in state y_1 (but may take any other values).
- “Virtual” or “likelihood” evidence: source of information is not sure about it.

Example:

- $e =$ Radiologist is 80% sure that $Xray=pos$
- Want e.g.:

$$P(Cancer|e) = P(Cancer|Xray, e)P(Xray|e) + P(Cancer|\neg Xray, e)P(\neg Xray|e)$$

- **Jeffrey Conditionalization**

Reasoning with numbers

- Reasoning with lung cancer example using Netica BN software.

(See Table 2.2 in *Bayesian AI* text.)

Understanding of Bayesian Networks (Semantics)

- A (more compact) representation of the joint probability distribution.
 - helpful in understanding how to construct network
- Encoding a collection of conditional independence statements.
 - helpful in understanding how to design inference procedures
 - via *Markov property / I-map*:
 - Each conditional independence implied by the graph is present in the probability distribution

Representing the joint probability distribution

Write $P(X_1 = x_1, \dots, X_n = x_n)$ as $P(x_1, x_2, \dots, x_n)$.

Factorization (chain rule):

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

Since BN structure implies that the value of a particular node is conditional *only* on the values of its parent nodes, this reduces to

$$P(x_1, x_2, \dots, x_n) = \prod_i P(x_i | \text{Parents}(X_i))$$

provided $\text{Parents}(X_i) \subseteq \{x_1, \dots, x_{i-1}\}$.

$$P(X = \text{pos} \wedge D = T \wedge C = T \wedge P = \text{low} \wedge S = F)$$

$$\begin{aligned} &= P(X = \text{pos} | D = T, C = T, P = \text{lo}, S = F) \\ &\quad \times P(D = T | C = T, P = \text{lo}, S = F) \\ &\quad \times P(C = T | P = \text{lo}, S = F) P(P = \text{lo} | S = F) P(S = F) \\ &= P(X = \text{pos} | C = T) P(D = T | C = T) \\ &\quad \times P(C = T | P = \text{lo}, S = F) P(P = \text{lo}) P(S = F) \end{aligned}$$

Pearl's Network Construction Algorithm

1. Choose the set of relevant variables $\{X_i\}$ that describe the domain.
2. Choose an ordering for the variables,
 $\langle X_1, \dots, X_n \rangle$.
3. While there are variables left:
 - (a) Add the next variable X_i to the network.
 - (b) Add arcs to the X_i node from some minimal set of nodes already in the net, $Parents(X_i)$, such that the following conditional independence property is satisfied:

$$P(X_i | X'_1, \dots, X'_m) = P(X_i | Parents(X_i))$$

where X'_1, \dots, X'_m are all the variables preceding X_i , including $Parents(X_i)$.

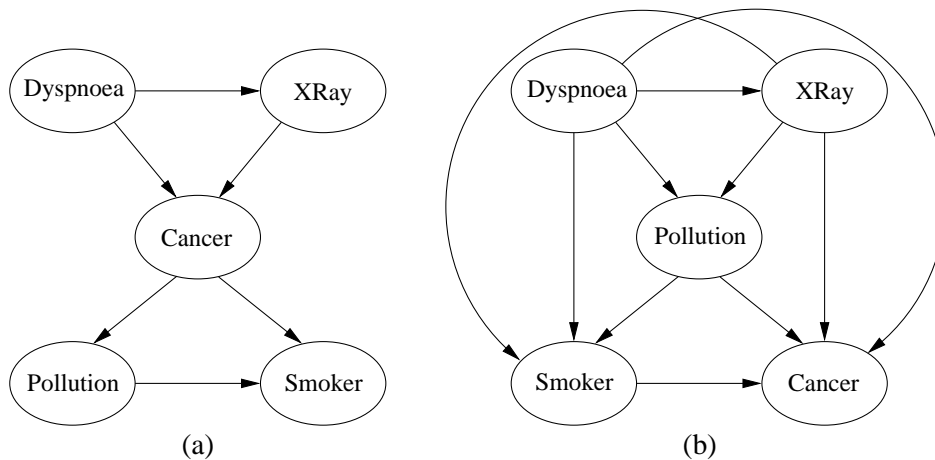
- (c) Define the CPT for X_i .

Compactness and Node Ordering

- Compactness of BN depends upon sparseness of the system.
- The best order to add nodes is to add the “root causes” first, then the variable they influence, so on until “leaves” reached.

→ Causal structure

- Alternative structures using different orderings (a) $\langle D, X, C, P, S \rangle$ (b) $\langle D, X, P, S, C \rangle$.



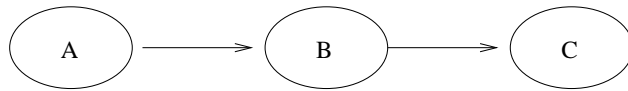
1. These BNs still represent same joint distribution.
2. Structure (b) requires as many probabilities as the full joint distribution! See below for *why*.

Conditional Independence

The relationship between conditional independence and BN structure is important for understanding how BNs work.

Conditional Independence: Causal Chains

Causal chains give rise to conditional independence:

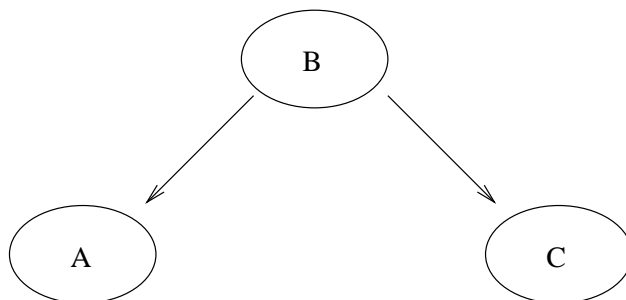


$$P(C|A \wedge B) = P(C|B)$$

Example: “smoking causes cancer which causes dyspnoea”

Conditional Independence: Common Causes

Common causes (or ancestors) also give rise to conditional independence:

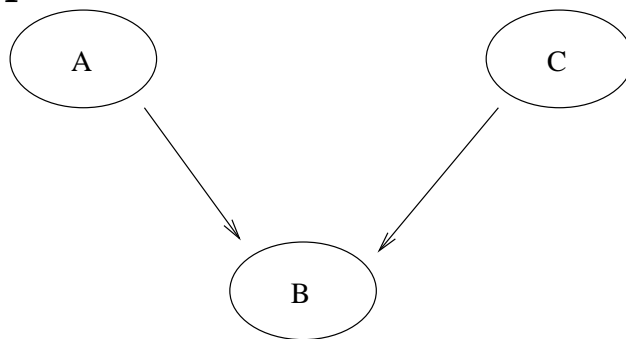


$$P(C|A \wedge B) = P(C|B) \equiv A \perp\!\!\!\perp C|B$$

Example: cancer is a common cause of the two symptoms, a positive XRay result and dyspnoea.

Conditional Dependence: Common Effects

Common effects (or their descendants) give rise to conditional *dependence*:



$$P(A|C \wedge B) \neq P(A)P(C) \equiv \neg(A \perp\!\!\!\perp C|B)$$

Example: Cancer is a common effect of pollution and smoking.

Given lung cancer, smoking “explains away” pollution.

D-separation

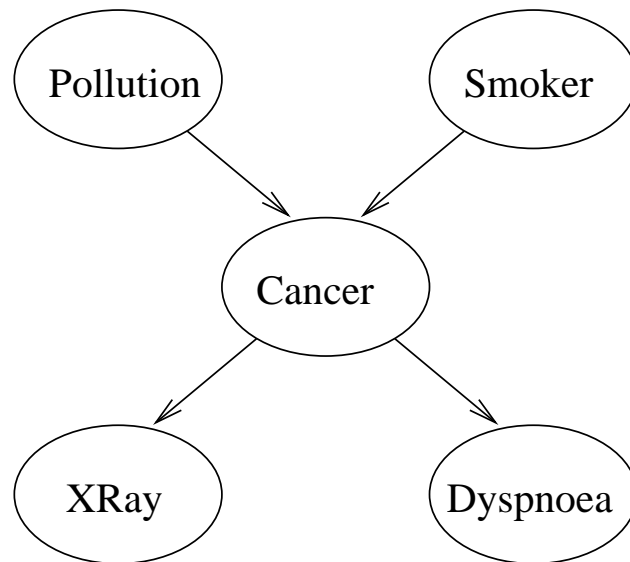
- Graphical criterion of conditional independence.

$$X \perp Y | Z$$

- We can determine whether a set of nodes X is independent of another set Y , given a set of evidence nodes E , via the Markov property:

$$X \perp Y | E \rightarrow X \perp\!\!\!\perp Y | E.$$

- Example



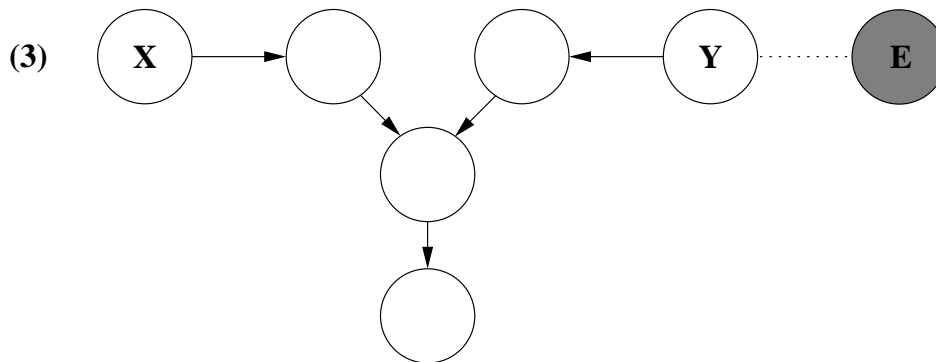
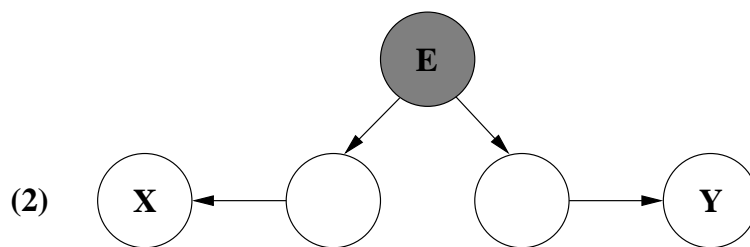
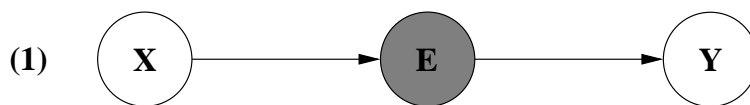
D-separation

How to determine d-separation, $X \perp Y | E$:

- If every undirected path from a node in X to a node in Y is *d-separated* by E , then X and Y are *conditionally independent* given E .
- A set of nodes E *d-separates* two sets of nodes X and Y if every undirected path from a node in X to a node in Y is *blocked* given E .
- A path is *blocked* given a set of nodes E if there is a node Z on the path for which one of three conditions holds:
 1. Z is in E and Z has one arrow on the path leading in and one arrow out (chain).
 2. Z is in E and Z has both path arrows leading out (common cause).
 3. Neither Z nor any descendant of Z is in E , and both path arrows lead in to Z (common effect).

D-separation (cont'd)

- Evidence nodes **E** shown shaded.



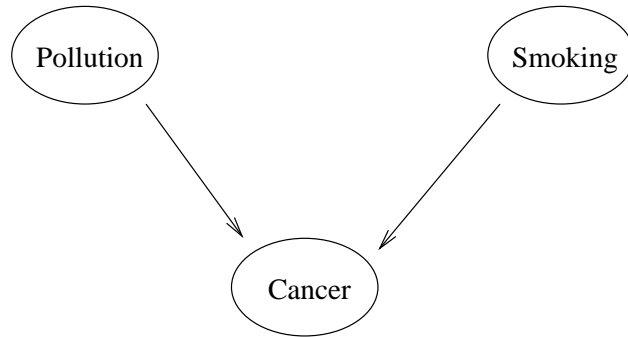
Causal Ordering

Why does variable order affect network density?

Because

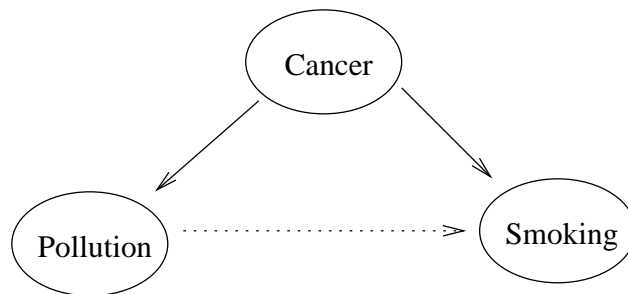
- Using the causal order allows direct representation of conditional independencies
- Violating causal order requires new arcs to re-establish conditional independencies

Causal Ordering (cont'd)



Pollution and *Smoking* are marginally independent.

Given the ordering: Cancer, Pollution, Smoking:



Marginal independence of *Pollution* and *Smoking* must be re-established by adding $Pollution \rightarrow Smoking$ or $Smoking \leftarrow Pollution$

Bayesian Networks: Summary

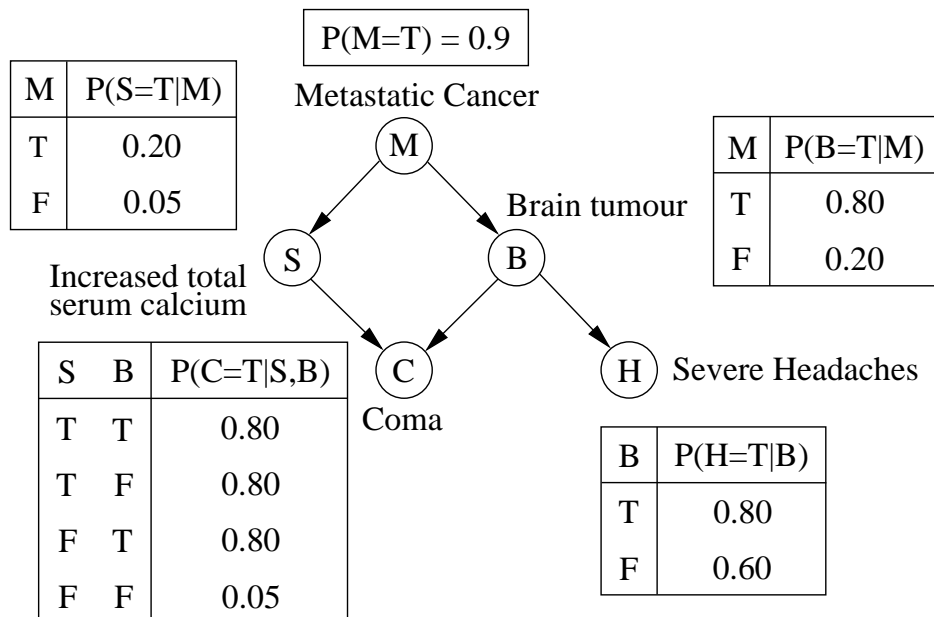
- Bayes' rule allows unknown probabilities to be computed from known ones.
- Conditional independence (due to causal relationships) allows efficient updating
- BNs are a natural way to represent conditional independence info.
 - links between nodes: qualitative aspects;
 - conditional probability tables: quantitative aspects.
- Probabilistic inference: compute the probability distribution for query variables, given evidence variables
- BN Inference is very flexible: can enter evidence about any node and update beliefs in any other nodes.

Inference Algorithms: Overview

- Exact inference
 - Trees and polytrees:
 - * message-passing algorithm
 - Multiply-connected networks:
- Approximate Inference
 - Large, complex networks:
 - * Stochastic Simulation
 - * Other approximation methods
- In the general case, both exact and approximate inference are computationally complex (“NP-hard”).
- Causal inference

Inference in multiply connected networks

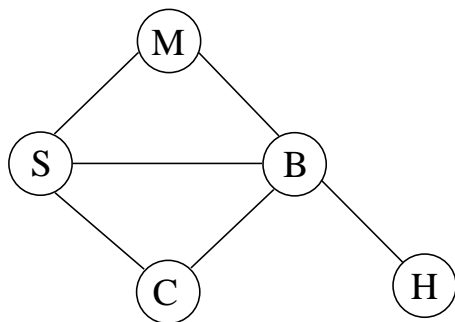
- Networks where two nodes are connected by more than one path
 - Two or more possible causes which share a common ancestor
 - One variable can influence another through more than one causal mechanism
- Example: Cancer network



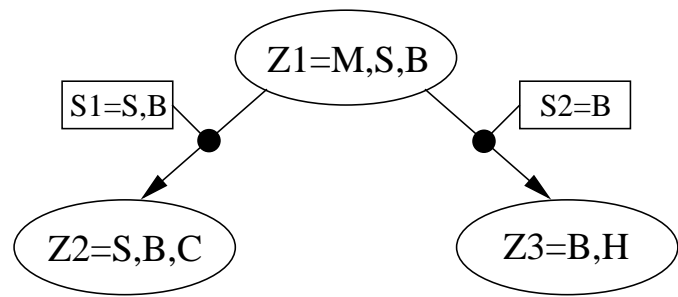
- Message passing doesn't work - evidence gets "counted twice"

Jensen join-tree method

- Jensen Join-tree (Jensen, 1996) version the current most efficient algorithm in this class (e.g. used in Hugin, Netica).



(a)



(b)

Jensen join-tree method (cont.)

Network evaluation done in two stages

- Compile into join-tree
 - May be slow
 - May require too much memory if original network is highly connected
- Do belief updating in join-tree (usually fast)

Caveat: clustered nodes have increased complexity; updates may be computationally complex

Causal modeling

We should like to do causal modeling with our Bayesian networks.

Prerequisite: arcs are truly causal (hence, nodes are properly ordered).

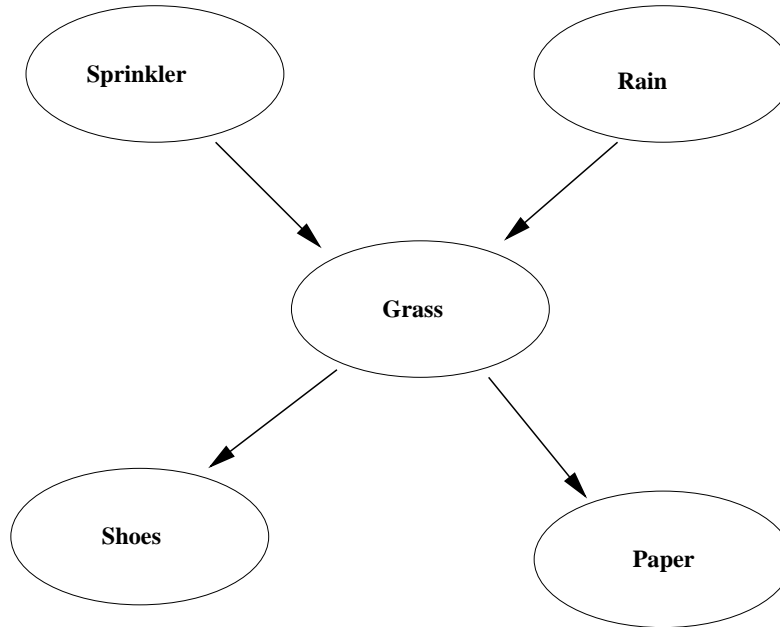
Reasoning about real or hypothetical interventions:

- what if we upgrade quality in manufacturing?
- what if we treat the patient with X, Y, Z?

For planning, control, prediction.

Common practice appears to be: let observation stand for intervention.

Causal inference



If we observe that the lawn is wet:

- We can infer in any direction; everything updates
- We get, e.g., “explaining away” between causes

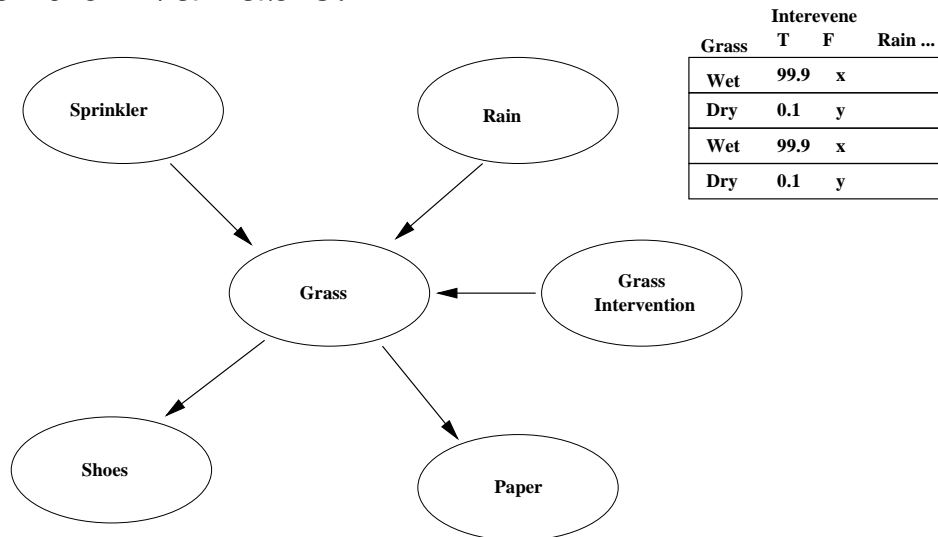
What happens if we *intervene* in a causal process?

Spirtes, et al., (1993), Pearl (2000) answer: cut links to parents and *then* update.

- No explaining away; parents are then unaffected
- Downstream updating is as normal

Causal inference

We prefer (conceptually) to augment the graph with an intervention variable:



- Simplistically, parent connections are severed
- With full generality, X acquires a new parent D_X
 - Allows any degree of control for intervention
 - Allows any kind of interaction with existing parents
 - Bayesian update algorithms unaffected

Inference: Summary

- Probabilistic inference: compute the probability distribution for query variables, given evidence variables
- Causal inference: compute the probability distribution for query variables, given intervention
- BN Inference is very flexible: can enter evidence about any node and update beliefs in any other nodes.
- The speed of inference in practice depends on the structure of the network: how many loops; numbers of parents; location of evidence and query nodes.
- BNs can be used to model causal intervention.

Extensions to Bayesian Networks

- For decision making: decision networks
- For reasoning about changes over time: dynamic Bayesian networks

Making Decisions

- Bayesian networks can be extended to support decision making.
- **Preferences** between different outcomes of various plans.
 - Utility theory
- **Decision theory** = Utility theory + Probability theory.

Expected Utility

$$EU(A|E) = \sum_i P(O_i|E, A) U(O_i|A) \quad (1)$$

- E = available evidence,
- A = a nondeterministic action
- O_i = possible outcome state
- U = utility

Decision Networks

A Decision network represents information about

- the agent's current state
- its possible actions
- the state that will result from the agent's action
- the utility of that state

Also called, *Influence Diagrams* (Howard&Matheson, 1981).

Type of Nodes

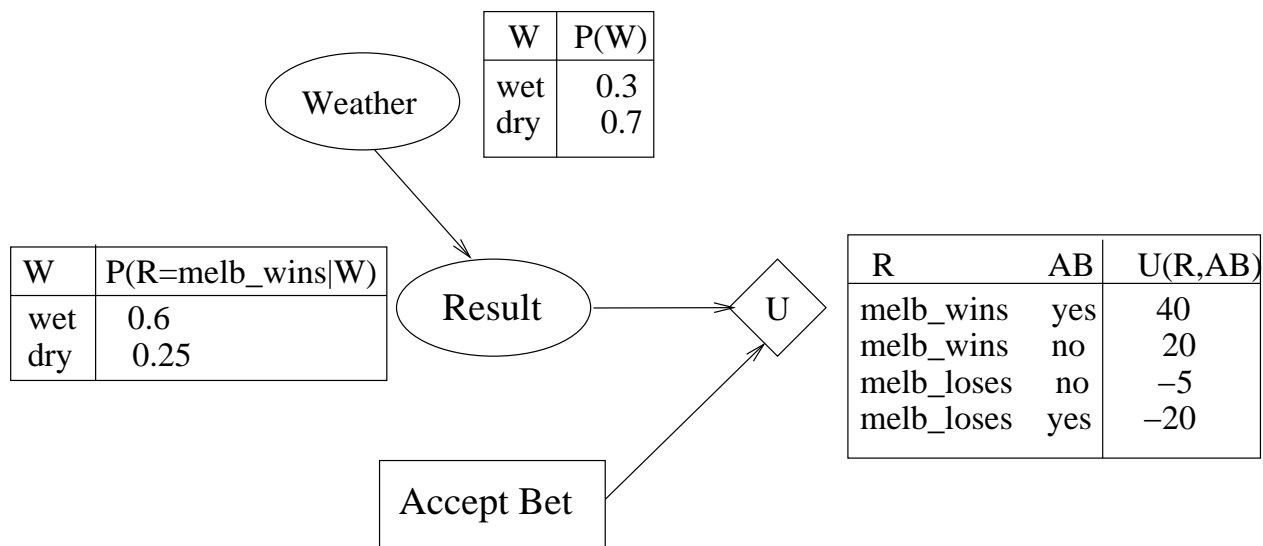
Chance nodes: (ovals) represent random variables (same as Bayesian networks). Has an associated CPT. Parents can be decision nodes and other chance nodes.

Decision nodes: (rectangles) represent points where the decision maker has a choice of actions.

Utility nodes: (diamonds) represent the agent's utility function (also called **value nodes** in the literature). Parents are variables describing the outcome state that directly affect utility. Has an associated table representing multi-attribute utility function.

Example: Football Team

Clare's football team, Melbourne, is going to play her friend John's team, Carlton. John offers Clare a friendly bet: whoever's team loses will buy the wine next time they go out for dinner. They never spend more than \$15 on wine when they eat out. When deciding whether to accept this bet, Clare will have to assess her team's chances of winning (which will vary according to the weather on the day). She also knows that she will be happy if her team wins and miserable if her team loses, regardless of the bet.



Evaluating Decision Networks: Algorithm

1. Add any available evidence.
2. For each action value in the decision node:
 - (a) Set the decision node to that value;
 - (b) Calculate the posterior probabilities for the parent nodes of the utility node, as for Bayesian networks, using a standard inference algorithm;
 - (c) Calculate the resulting expected utility for the action.
3. Return the action with the highest expected utility.

Evaluating Decision Networks: Example

$$\begin{aligned}P(R = melb_wins) &= P(W = w)P(R = melb_wins|W = w) \\ &= +P(W = d)P(R = melb_wins|W = d)\end{aligned}$$

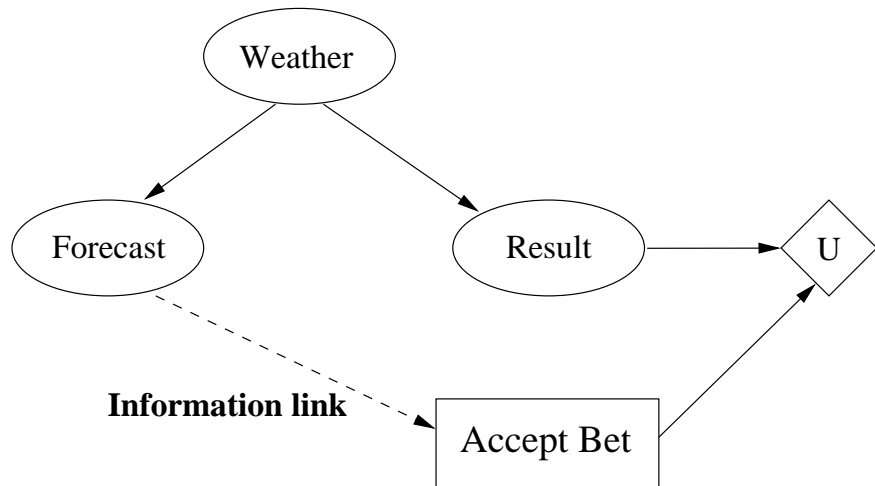
$$\begin{aligned}EU(AB = yes) &= P(R = wins)U(R = wins|AB = yes) \\ &+ P(R = loses)U(R = loses|AB = yes) \\ &= (0.3 \times 0.6 + 0.7 \times 0.25)40 \\ &+ (0.3 \times 0.4 + 0.7 \times 0.75) - 20 \\ &= 0.355 \times 40 + 0.645 \times -20 = 14.2 - 12.9 \\ &= 1.3\end{aligned}$$

$$\begin{aligned}EU(AB = no) &= P(R = wins)U(R = wins|AB = no) \\ &+ P(R = loses)U(R = loses|AB = no) \\ &= (0.3 \times 0.6 + 0.7 \times 0.25)20 \\ &+ (0.3 \times 0.4 + 0.7 \times 0.75) - 5 \\ &= 0.355 \times 20 + 0.645 \times -5 = 7.1 - 3.225 \\ &= 3.875\end{aligned}$$

Information Links

- Indicate when a chance node needs to be observed before a decision is made.

W	F	P(F W)
wet	rainy	0.60
	cloudy	0.25
	sunny	0.15
dry	rainy	0.10
	cloudy	0.40
	sunny	0.50



Decision Table

F	Accept Bet
rainy	yes
cloudy	no
sunny	no

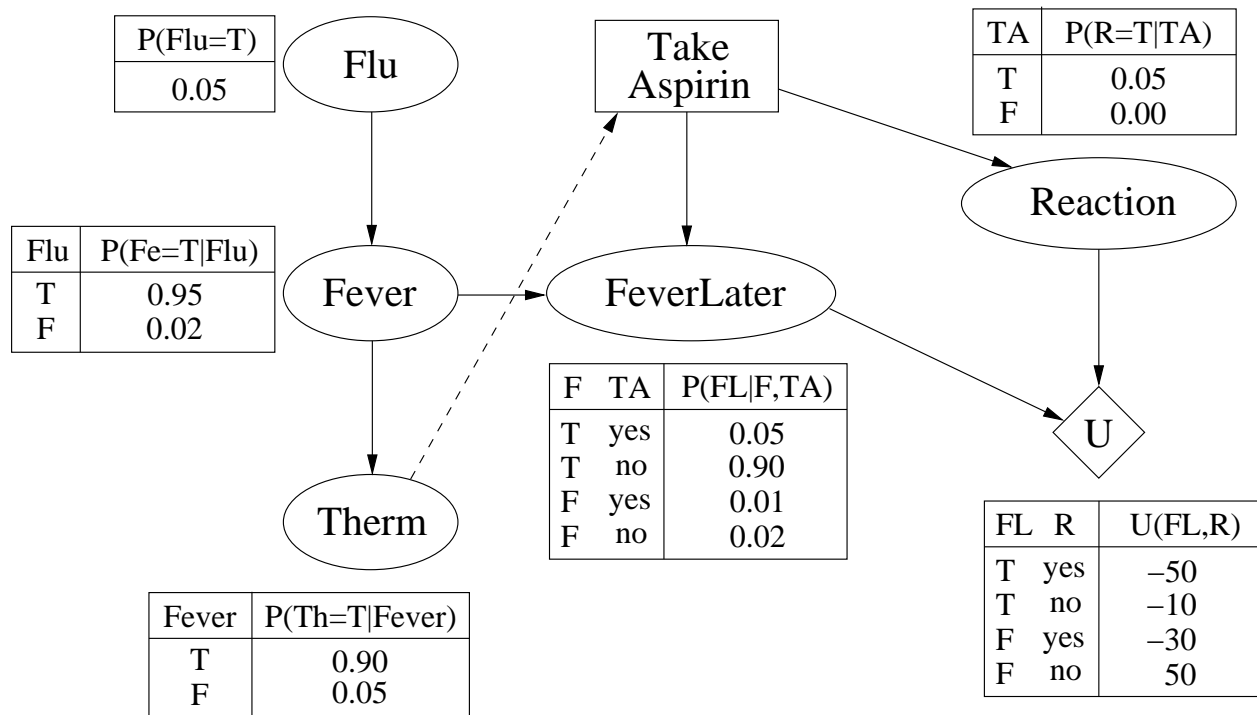
Decision Table Algorithm

1. Add any available evidence.
2. For each combination of values of the parents of the decision node:
 - (a) For each action value in the decision node:
 - i. Set the decision node to that value;
 - ii. Calculate the posterior probabilities for the parent nodes of the utility node, as for Bayesian networks, using a standard inference algorithm;
 - iii. Calculate the resulting expected utility for the action.
 - (b) Record the action with the highest expected utility in the decision table.
3. Return the decision table.

Fever problem description

Suppose that you know that a fever can be caused by the flu. You can use a thermometer, which is fairly reliable, to test whether or not you have a fever. Suppose you also know that if you take aspirin it will almost certainly lower a fever to normal. Some people (about 5% of the population) have a negative reaction to aspirin. You'll be happy to get rid of your fever, as long as you don't suffer an adverse reaction if you take aspirin.

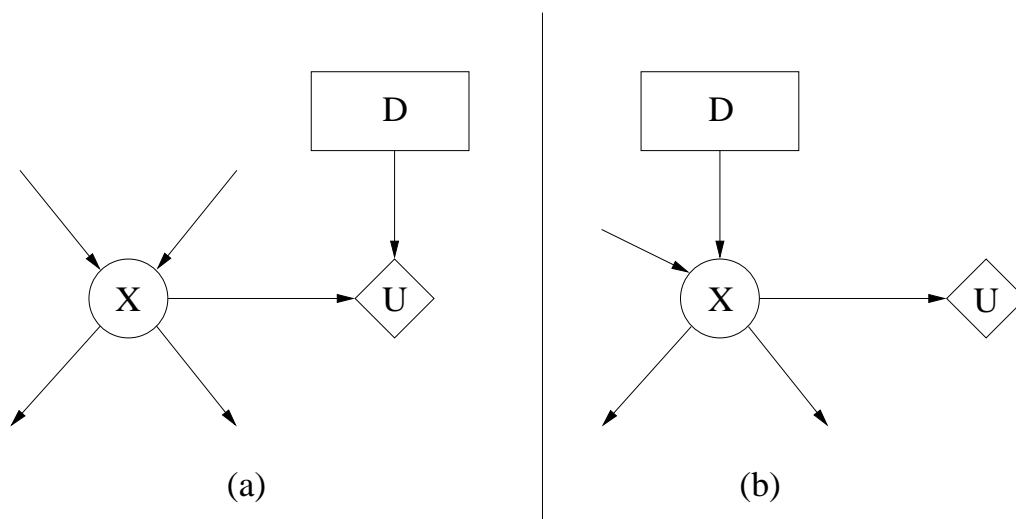
Fever decision network



Fever decision table

Ev.	$Bel(FL=T)$	EU(TA=yes)	EU(TA=no)	Dec.
None	0.046	45.27	45.29	no
$Th=F$	0.525	45.41	48.41	no
$Th=T$	0.273	44.1	19.13	yes
$Th=T$ & $Re=T$	0.273	-30.32	0	no

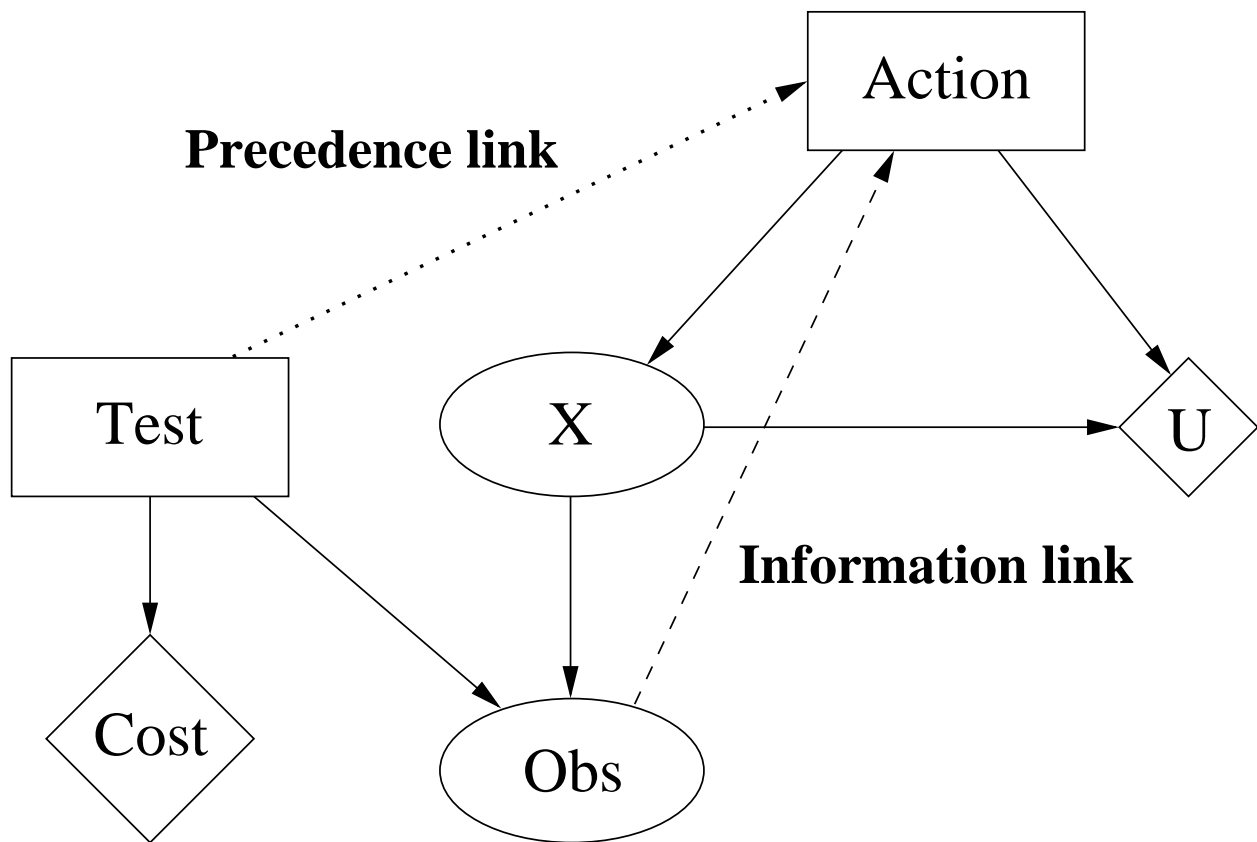
Types of actions



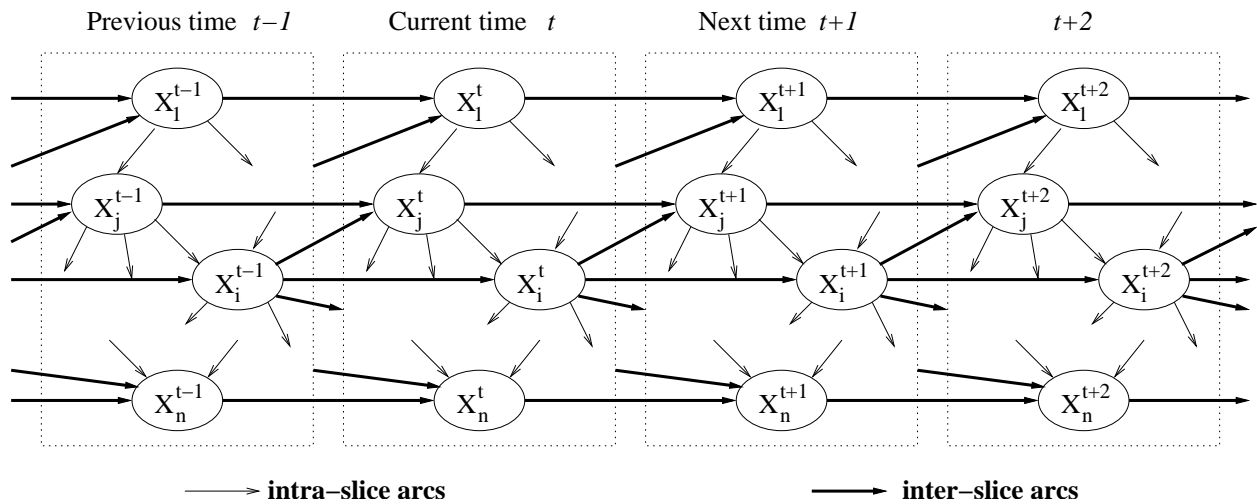
(a) Non-intervening and (b) Intervening

Sequential decision making

- Precedence links used to show temporal ordering.
- Network for a test-action decision sequence

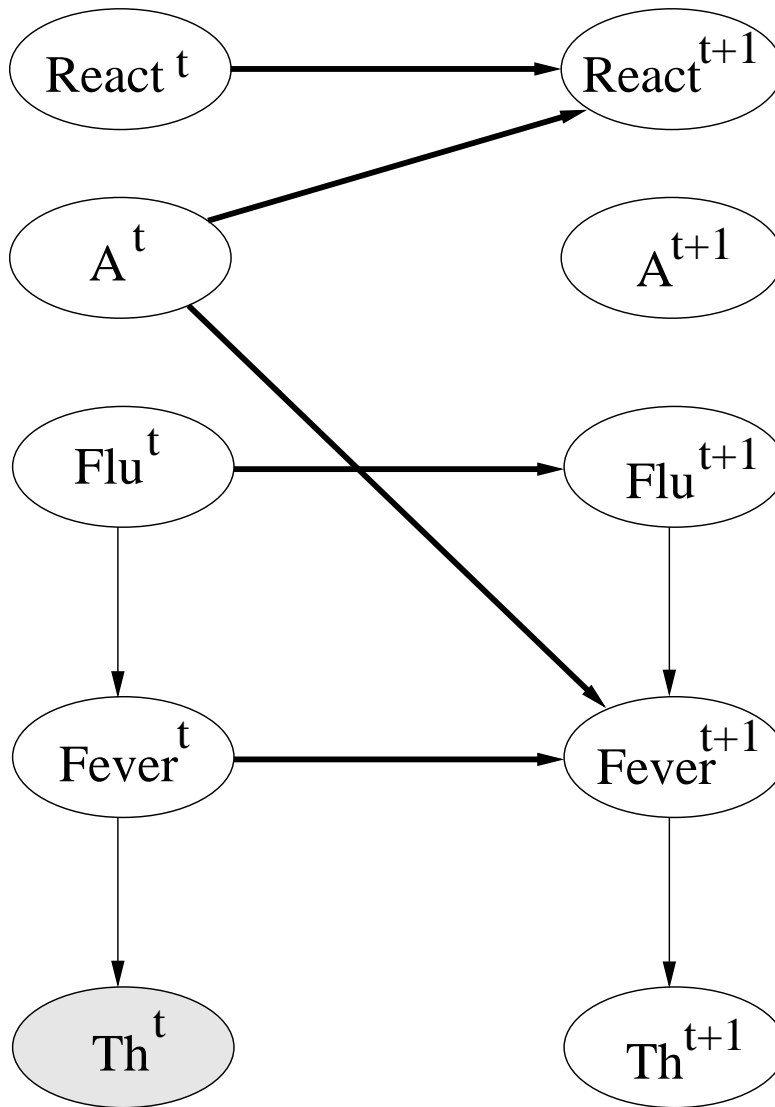


Dynamic Belief Networks



- One node for each variable for each time step.
- **Intra-slice arcs** $X_i^T \longrightarrow X_j^T$
- **Inter-slice (temporal) arcs**
 1. $X_i^T \longrightarrow X_i^{T+1}$
 2. $X_i^T \longrightarrow X_j^{T+1}$

Fever DBN

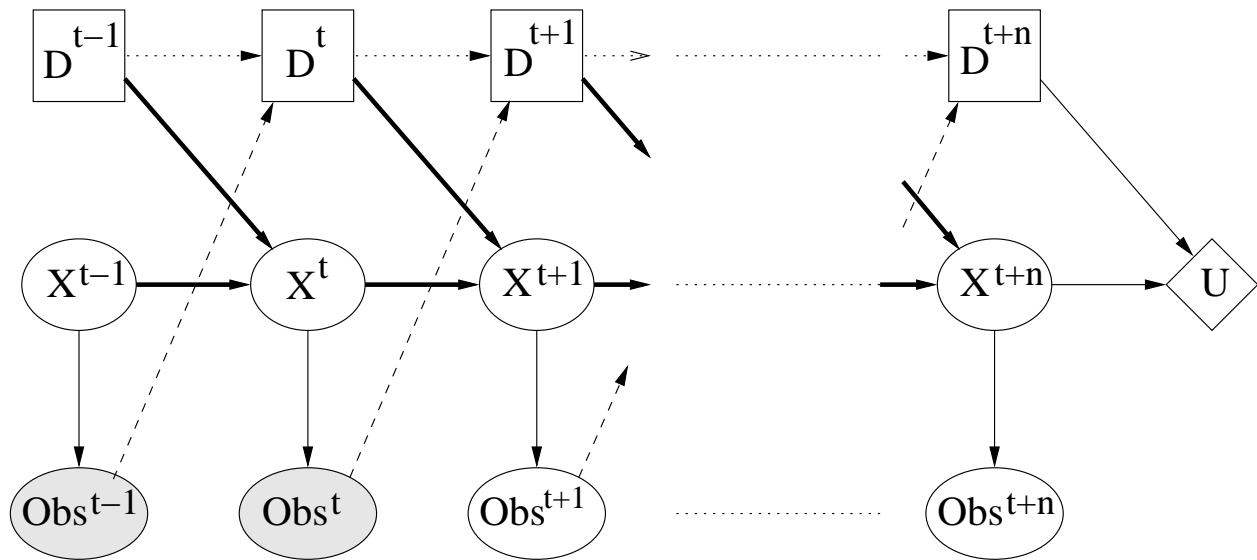


DBN reasoning

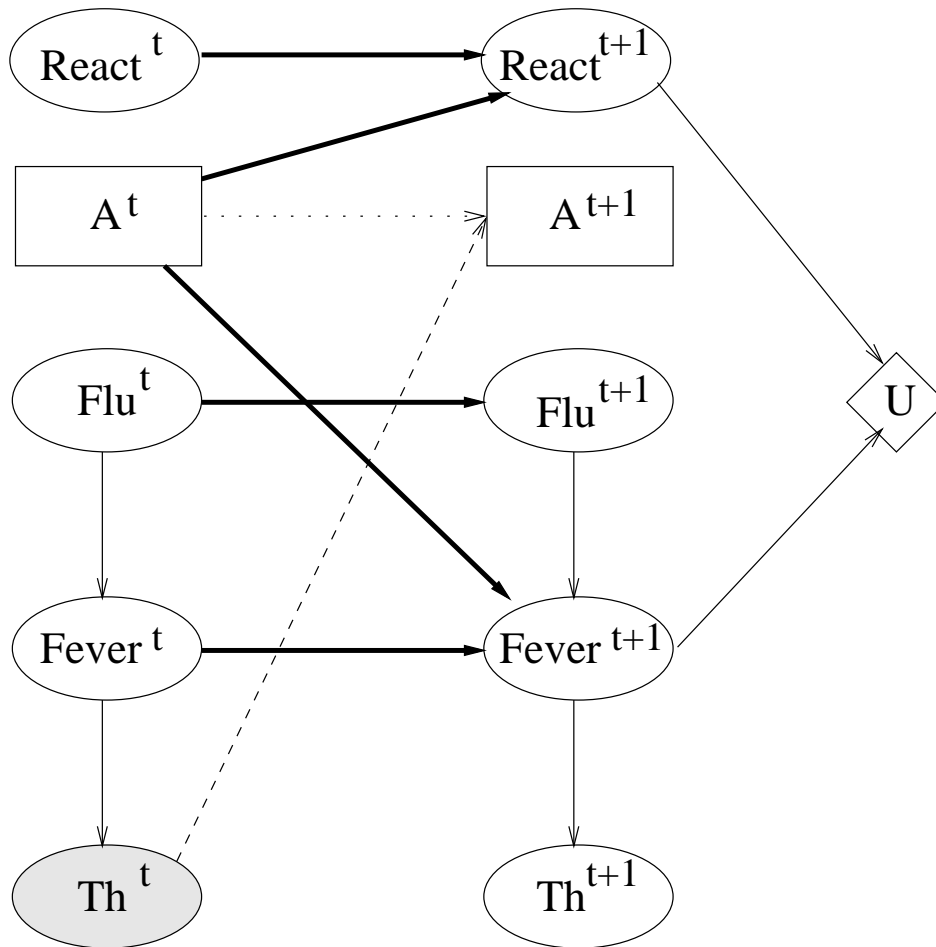
- Can calculate distributions for S_{t+1} and further: **probabilistic projection.**
- Reasoning can be done using standard BN updating algorithms
- This type of DBN gets very large, very quickly.
- Usually only keep two time slices of the network.

Dynamic Decision Network

- Similarly, Decision Networks can be extended to include temporal aspects.
- Sequence of decisions taken = Plan.



Fever DDN



Extensions: Summary

- BNs can be extended with decision nodes and utility nodes to support decision making: *Decision Networks* or *Influence Diagrams*.
- BNs and decision networks can be extended to allow explicit reasoning about changes over time.

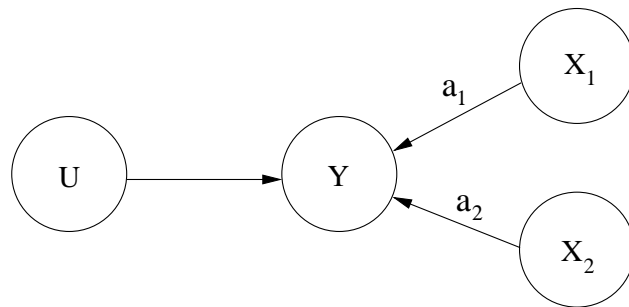
Learning Bayesian Networks

- Linear and Discrete Models
- Learning Network Parameters
 - Linear Coefficients
 - Learning Probability Tables
- Learning Causal Structure
- Conditional Independence Learning
 - Statistical Equivalence
 - TETRAD II
- Bayesian Learning of Bayesian Networks
 - Cooper & Herskovits: K2
 - Learning Variable Order
 - Statistical Equivalence Learners
- Full Causal Learners
- Minimum Encoding Methods
 - Lam & Bacchus's MDL learner
 - MML metrics
 - MML search algorithms
 - MML Sampling
- Empirical Results

Linear and Discrete Models

Linear Models: Used in biology & social sciences since Sewall Wright (1921)

Linear models represent causal relationships as sets of linear functions of “independent” variables.



Equivalently:

$$X_3 = a_{13}X_1 + a_{23}X_2 + \epsilon_1$$

Structural equation models (SEMs) are close relatives

Discrete models: “Bayesian nets” replace vectors of linear coefficients with CPTs.

Learning Causal Structure

This is the *real* problem; parameterizing models is relatively straightforward estimation problem.

Size of the dag space is superexponential:

- Number of possible orderings: $n!$
- Times number of ways of pairing up (for arcs): $2^{C_2^n}$
- Minus number of possible cyclic graphs

Without the subtraction (which is a small proportion):

n	$n!2^{C_2^n}$
0	0
1	1
2	4
3	48
4	1536
5	12280
10	127677049435953561600
100	[too many digits to show]

Learning Causal Structure

There are two basic methods:

- Learning from conditional independencies (CI learning)
- Learning using a scoring metric (Metric learning)

CI learning (Verma and Pearl, 1991)

Suppose you have an Oracle who can answer yes or no to any question of the type:

$$X \perp\!\!\!\perp Y | S?$$

(i.e., is X conditional independence Y given S)

Then you can learn the correct causal model, up to statistical equivalence (patterns).

Verma-Pearl Algorithm

two rules allow discovery of the set of causal models consistent with all such answers (“patterns”):

1. **Principle I** Put an undirected link between any two variables X and Y iff for every \mathbf{S} s.t. $X, Y \notin \mathbf{S}$

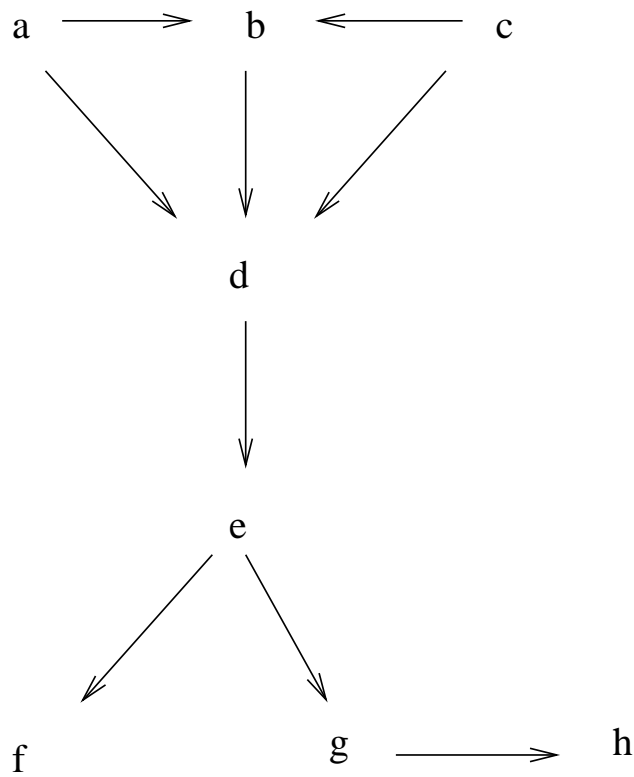
$$\neg(X \perp\!\!\!\perp Y) | \mathbf{S}$$

2. **Principle II** For every undirected v-structure $X - Z - Y$ orient the arcs $X \rightarrow Z \leftarrow Y$ iff

$$\neg(X \perp\!\!\!\perp Y) | \mathbf{S}$$

for **every** \mathbf{S} s.t. $X, Y \notin \mathbf{S}$ and $Z \in \mathbf{S}$.

CI learning example



CI learning example

$$1) a - b - c \quad a \rightarrow b \leftarrow c$$

b [induces a dependency]

$$2) a - d - c \quad a \rightarrow d \leftarrow c$$

$$3) c - d - e \quad \neg(c \rightarrow d \leftarrow e)$$

therefore $c \rightarrow d \rightarrow e$

$$4) a - d - e \quad \text{no news}$$

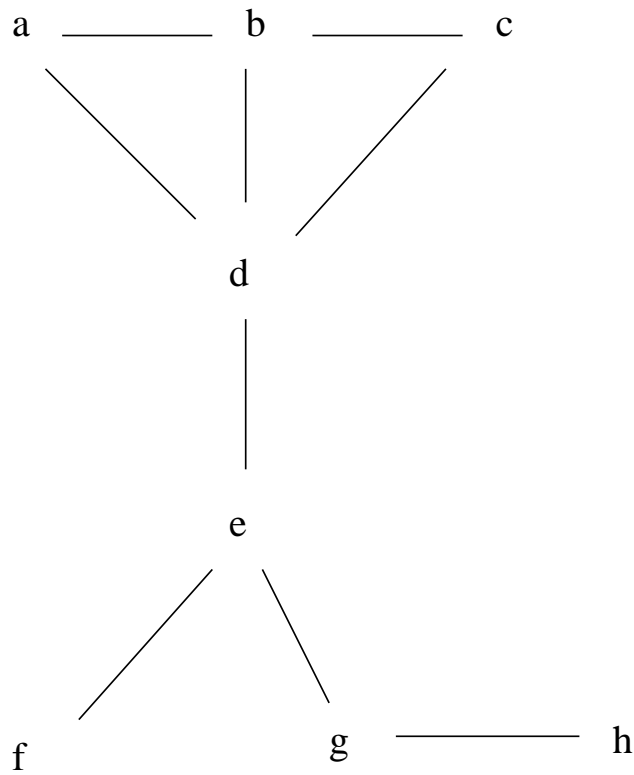
$$5) b - d - e \quad \text{no news}$$

$$6) d - e - f \quad \neg(d \rightarrow e \leftarrow f)$$

$$7) d - e - g \quad \neg(d \rightarrow e \leftarrow g)$$

$$6) e - g - h \quad \neg(e \rightarrow g \leftarrow h)$$

CI learning example



Statistical Equivalence

Verma and Pearl's rules identify the set of causal models which are statistically equivalent —

Two causal models H_1 and H_2 are **statistically equivalent** iff they contain the same variables and joint samples over them provide no statistical grounds for preferring one over the other.

Examples

- All fully connected models are equivalent.
- $A \rightarrow B \rightarrow C$ and $A \leftarrow B \leftarrow C$.
- $A \rightarrow B \rightarrow D \leftarrow C$ and $A \leftarrow B \rightarrow D \leftarrow C$.

Statistical Equivalence

- (Verma and Pearl, 1991): Any two causal models over the same variables which have the same skeleton (undirected arcs) and the same directed v-structures are statistically equivalent.
- Chickering (1995): If H_1 and H_2 are statistically equivalent, then they have the same maximum likelihoods relative to any joint samples:

$$\max P(e|H_1, \theta_1) = \max P(e|H_2, \theta_2)$$

where θ_i is a parameterization of H_i

TETRAD II

— Spirtes, Glymour and Scheines (1993)

Replace the Oracle with statistical tests:

- for linear models a significance test on partial correlation

$$X \perp\!\!\!\perp Y | \mathbf{S} \text{ iff } \rho_{XY \cdot \mathbf{S}} = 0$$

- for discrete models a χ^2 test on the difference between CPT counts expected with independence (E_i) and observed (O_i)

$$X \perp\!\!\!\perp Y | \mathbf{S} \text{ iff } \sum_i O_i \ln \left(\frac{O_i}{E_i} \right)^2 \approx 0$$

Implemented in their **PC Algorithm**

Bayesian LBN: Cooper & Herskovits' K2

— Cooper & Herskovits (1991, 1992)

Compute $P(h_i|e)$ by brute force, under the assumptions:

1. All variables are discrete.
2. Samples are i.i.d.
3. No missing values.
4. All values of child variables are uniformly distributed.
5. Priors over hypotheses are uniform.

With these assumptions, Cooper & Herskovits reduce the computation of $P_{CH}(h, e)$ to a polynomial time counting problem.

Cooper & Herskovits

But the hypothesis space is exponential; they go for dramatic simplification:

6. Assume we know the temporal ordering of the variables.

In that case, for any pair of variables the only problem is

- deciding whether they are connected by an arc
 - arc direction is trivial
 - cycles are impossible.

New hypothesis space has size only $2^{(n^2-n)/2}$ (still exponential).

Algorithm “K2” does a greedy search through this reduced space.

Learning Variable Order

Reliance upon a given variable order is a major drawback to K2

And many other algorithms (Buntine 1991, Bouckert 1994, Suzuki 1996, Madigan & Raftery 1994)

What's wrong with that?

- We want autonomous AI (data mining). If experts can order the variables they can likely supply models.
- Determining variable ordering is half the problem. If we know A comes before B , the only remaining issue is whether there is a link between the two.
- The number of orderings consistent with dags is exponential (Brightwell & Winkler 1990; number complete). So iterating over all possible orderings will not scale up.

Statistical Equivalence Learners

Heckerman & Geiger (1995) advocate learning only up to statistical equivalence classes (a la TETRAD II).

Since observational data cannot distinguish btw equivalent models, there's no point trying to go further.

⇒ Madigan, Andersson, Perlman & Volinsky (1996) follow this advice, use uniform prior over equivalence classes.

⇒ Geiger and Heckerman (1994) define Bayesian metrics for linear and discrete equivalence classes of models (BGe and BDe)

GES

Greedy Equivalence Search (GES)

- Product of the CMU-Microsoft group (Meek, 1996; Chickering, 2002)
- Two-stage greedy search: Begin with unconnected pattern
 1. Greedily add single arcs until reaching a local maximum
 2. Prune back edges which don't contribute to the score
- Uses a Bayesian score over patterns only
- Implemented in TETRAD and Murphy's BNT

Statistical Equivalence Learners

Wallace & Korb (1999): This is not right!

- These are **causal** models; they *are* distinguishable on *experimental* data.
 - Failure to collect some data is no reason to change prior probabilities.
E.g., If your thermometer topped out at 35° , you wouldn't treat $\geq 35^\circ$ and 34° as equally likely.
- Not all equivalence classes are created equal:
 $\{ A \leftarrow B \rightarrow C, A \rightarrow B \rightarrow C, A \leftarrow B \leftarrow C \}$
 $\{ A \rightarrow B \leftarrow C \}$
- *Within* classes some dags should have greater priors than others... E.g.,
LightsOn \rightarrow InOffice \rightarrow LoggedOn v.
LightsOn \leftarrow InOffice \rightarrow LoggedOn

Full Causal Learners

So... a full causal learner is an algorithm that:

1. Learns causal connectedness.
 2. Learns v-structures.
Hence, learns equivalence classes.
 3. Learns full variable order.
Hence, learns full causal structure (order + connectedness).
- TETRAD II: 1, 2.
 - Madigan et al.; Heckerman & Geiger (BGe, BDe): 1, 2.
 - GES: 1, 2.
 - Cooper & Herskovits' K2: 1.
 - Lam and Bacchus MDL: 1, 2 (partial), 3 (partial).
 - Wallace, Neil, Korb MML: 1, 2, 3.

CaMML

Minimum Message Length (Wallace & Boulton 1968)
uses Shannon's measure of information:

$$I(m) = -\log P(m)$$

Applied in reverse, we can compute $P(h, e)$ from $I(h, e)$.

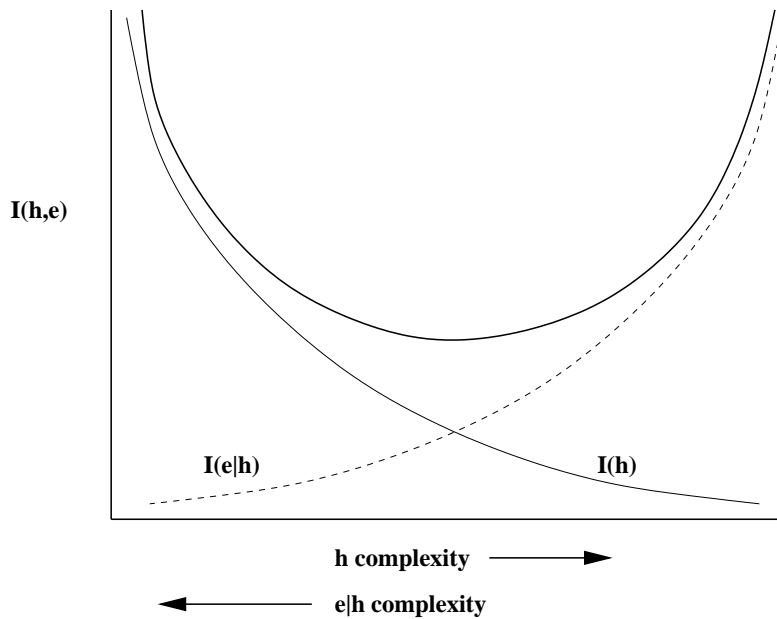
Given an *efficient* joint encoding method for the hypothesis & evidence space (i.e., satisfying Shannon's law), MML:

Searches $\{h_i\}$ for that hypothesis h that
minimizes $I(h) + I(e|h)$.

Applies a trade-off between

- Model simplicity
- Data fit

MML Metric



Equivalent to that h that maximizes $P(h)P(e|h)$ — i.e., $P(h|e)$.

$$\begin{aligned} I(h, e) &= I(h) + I(e|h) \\ -\log P(h, e) &= -\log P(h) - \log P(e|h) \\ -\log P(h, e) &= -\log P(h)P(e|h) \\ P(h, e) &= P(h)P(e|h) \end{aligned}$$

Hence, $\min I(h, e) \equiv \max P(h, e)$.

MML Metric for discrete models

We can use $P_{CH}(h_i, e)$ (from Cooper & Herskovits) to define an MML metric for discrete models.

Difference between MML and Bayesian metrics:

MML partitions the parameter space and selects optimal parameters.

Equivalent to a penalty of $\frac{1}{2} \log \frac{\pi e}{6}$ per parameter (Wallace & Freeman 1987); hence:

$$I(e, h_i) = \frac{s_j}{2} \log \frac{\pi e}{6} - \log P_{CH}(h_i, e) \quad (2)$$

Applied in MML Sampling algorithm.

MML Sampling

Search space of totally ordered models (TOMs).

Sampled via a Metropolis algorithm (Metropolis et al. 1953).

From current model M , find the next model M' by:

- Randomly select a variable; attempt to swap order with its predecessor.
- Or, randomly select a pair; attempt to add/delete an arc.

Attempts succeed whenever $P(M')/P(M) > U$ (per MML metric), where U is uniformly random from $[0 : 1]$.

MML Sampling

Metropolis: this procedure samples TOMs with a frequency proportional to their posterior probability.

To find posterior of dag h : keep count of visits to all TOMs consistent with h

Estimated by counting visits to all TOMs with identical max likelihoods to h

Output: Probabilities of

- Top dags
- Top statistical equivalence classes
- Top MML equivalence classes

KEBN: Overview

- The BN Knowledge Engineering Process
- Model construction
 - Variables and values
 - Graph Structure
 - Probabilities
 - Preferences
- Evaluation

Knowledge Engineering with Bayesian Networks (KEBN)

(Laskey, 1999).

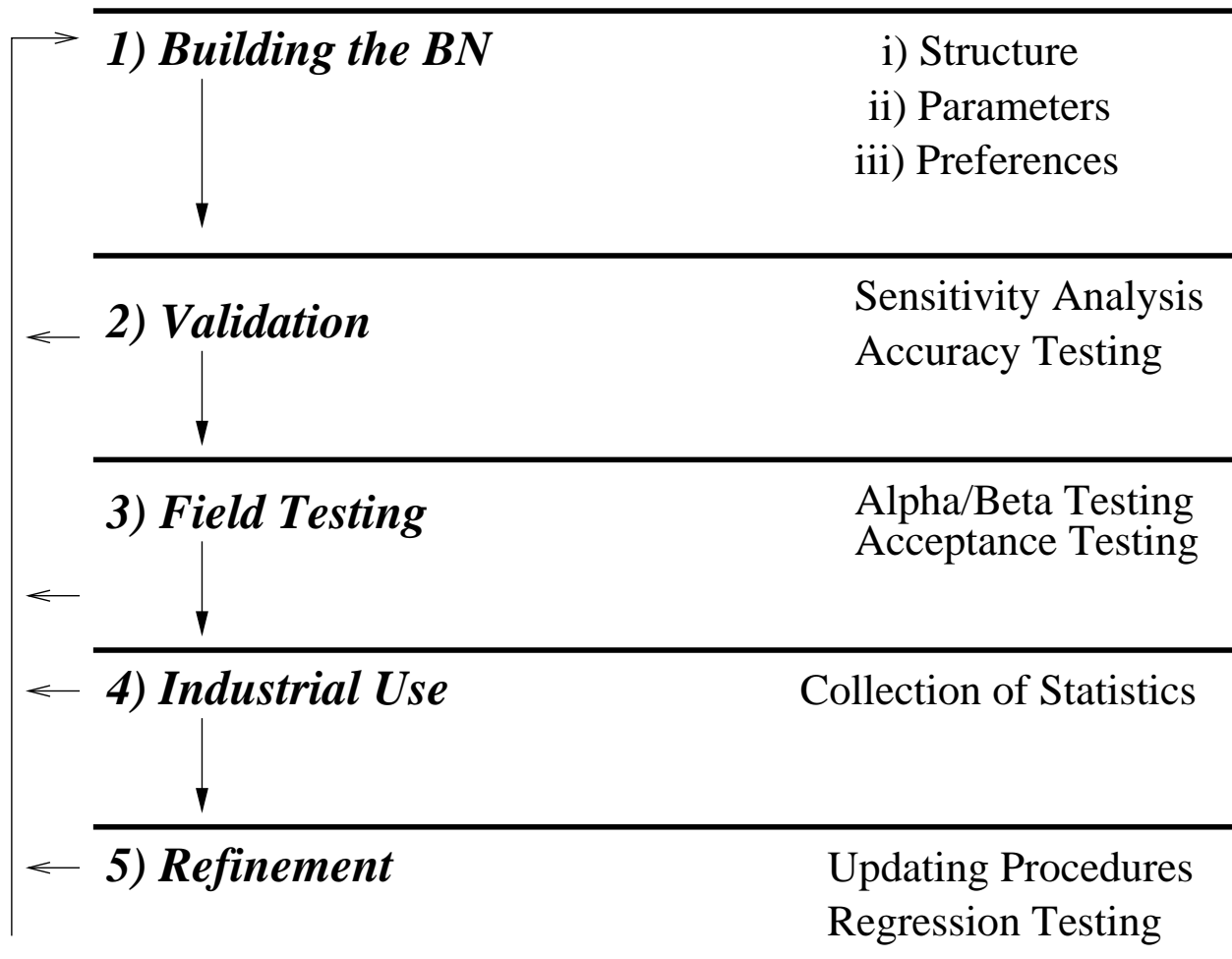
- Objective: Construct a model to perform a defined task
- Participants: Collaboration between domain expert(s) and BN modelling expert(s), including use of automated methods.
- Process: iterate until “done”
 - Define task objective
 - Construct model
 - Evaluate model

KEBN

Production of Bayesian/decision nets for

- **Decision making:** Which policy carries the least risk of failure?
- **Forward Prediction:** Hypothetical or factual. Who will win the election?
- **Retrodiction/Diagnosis:** Which illness do these symptoms indicate?
- **Monitoring/control:** Do containment rods need to be inserted here at Chernobal?
- **Explanation:** Why did the patient die? Which cause exerts the greater influence?
- **Sensitivity Analysis:** What range of probs/utilities make no difference to X?
- **Information value:** What's the differential utility for changing precision of X to ϵ ?

KEBN Lifecycle Model



Notes on Lifecycle Model

- Phase 1: Building Bayesian Networks.
 - Major network components: structure, parameters and utilities.
 - Elicitation: from experts, learned with data mining methods, or some combination of the two.
- Phase 2: Evaluation.
 - Networks need to be validated for: predictive accuracy, respecting known temporal order of the variables and respecting known causal structure.
 - Use statistical data (if available) or expert judgement.
- Phase 3: Field Testing.
 - Domain expert use BN to test usability, performance, etc.

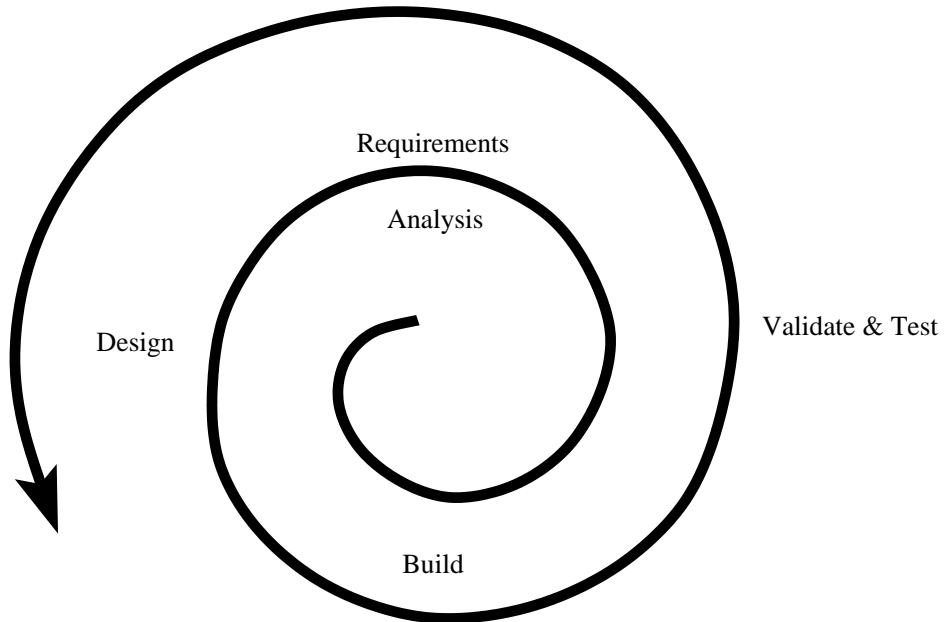
Notes on Lifecycle Model (cont.)

- Phase 4: Industrial Use.
 - Requires a statistics collection regime for on-going validation and/or refinement of the networks.
- Phase 5: Refinement.
 - Requires a process for receiving and incorporating change i requests
 - Includes regression testing to verify that changes do not undermine established performance.

KEBN Spiral Model

From Laskey & Mahoney (2000)

Idea (from Boehm, Brooks): prototype-test cycle



KEBN Tasks

For Bayesian Networks, identifying:

1. What are the variables? What are their values/states?
2. What is the graph structure? What are the direct causal relationships?
3. What are the parameters (probabilities)? Is there local model structure?

When building decision nets, identifying:

4. What are the available actions/decisions?
5. What are the utility nodes & their dependencies?
6. What are the preferences (utilities)?

The major methods are:

- Expert elicitation (1-6)
- Automated learning from data (1-3, 5-6?)
- Adapting from data (1-3, 5-6?)

Variables

Which are the most important variables?

- “Focus” or “query” variables
 - variables of interest
- “Evidence” or “observation” variables
 - What sources of evidence are available?
- “Context” variables
 - Sensing conditions, background causal conditions
- “Controllable” variables
 - variables that can be “set”, by intervention

Start with query variables and spread out to related variables.

NB: Roles of variables may change.

Variable values/states

- Variable values must be exclusive and exhaustive
 - Naive modelers sometimes create separate (often Boolean) variables for different states of the same variable
- Types of variables
 - Binary (2-valued, including Boolean)
 - Qualitative
 - Numeric discrete
 - Numeric continuous
- Dealing with infinite and continuous variable domains
 - Some BN software (e.g. Netica) requires that continuous variables be discretized
 - Discretization should be based on differences in effect on related variables (i.e. not just be even sized chunks)

Graphical structure

Goals in specifying graph structure

- Minimize probability elicitation: fewer nodes, fewer arcs, smaller state spaces
- Maximize fidelity of model
 - Sometimes requires more nodes, arcs, states
 - Tradeoff between more accurate model and cost of additional modelling
 - Too much detail can decrease accuracy
- Drawing arcs in causal direction is not “required” BUT
 - Increases conditional independence
 - Results in more compact model
 - Improves ease of probability elicitation
- If mixing continuous and discrete variables
 - Exact inference algorithms only for the case where discrete variables are ancestors, not descendants of continuous variables

Relationships between variables

Types of qualitative understanding can help determine local/global structure

- Causal relationships
 - Variables that could cause a variable to take a particular state
 - Variables that could prevent a variable taking a particular state
- Enabling variables
 - Conditions that permit, enhance or inhibit operation of a cause
- Effects of a variable
- Associated variables
 - When does knowing a value provide information about another variable?

Relationships between variables (cont.)

- Dependent and independent variables
 - D-separation tests
 - Which pairs are directly connected?
Probabilities dependent regardless of all other variables?

Matilda - software tool for visual exploration of dependencies (Boneh, 2002)

- Temporal ordering of variables
- Explaining away/undermining
- Causal non-interaction/additivity
- Causal interaction
 - Positive/negative Synergy
 - Preemption
 - Interference/XOR
- Screening off: causal proximity
- Explanatory value
- Predictive value

Probabilities

- The parameters for a BN are a set of conditional probability distributions of child values given values of parents
- One distribution for each combination of values of parent variables
- Assessment is exponential in the number of parent variables
- The number of parameters can be reduced by taking advantage of additional structure in the domain (called **local model structure**)

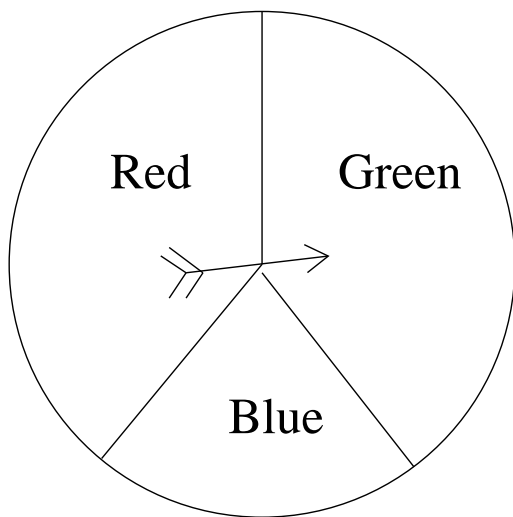
Probability Elicitation

- Discrete variables
 - Direct elicitation: “ $p=0.7$ ”
 - Odds (esp. for very small probs): “1 in 10,000”
 - Qualitative assessment: “very high probability”
 - * Use scale with numerical and verbal anchors (van der Gaag et al., 1999)
 - * Do mapping separately from qualitative assessment
- Continuous variables
 - bi-section method
 - * Elicit median: equally likely to be above and below
 - * Elicit 25th percentile: bisects interval below median
 - * Continue with other percentiles till fine enough discriminations
 - Often useful to fit standard functional form to expert’s judgements
 - Need to discretize for most BN software

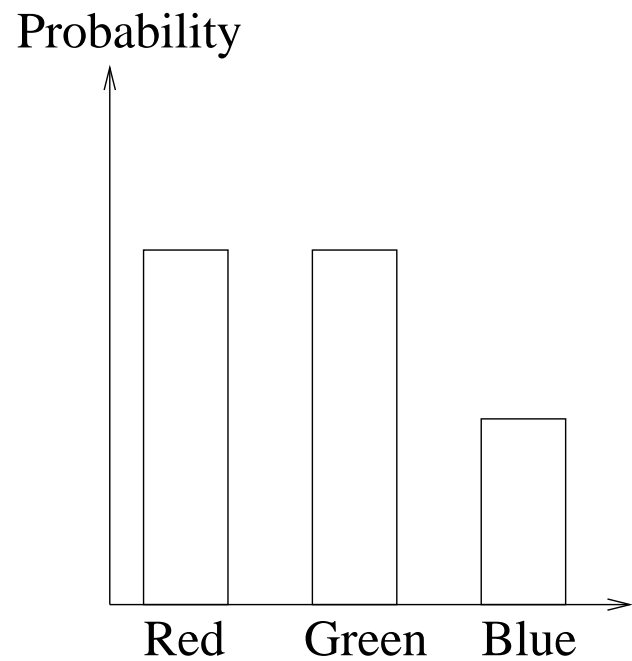
Probability Elicitation

Graphical aids are known to be helpful

- pie charts
- histograms



(a)



(b)

Probability Elicitation (cont.)

- Combination of qualitative and quantitative assessment
- Automated correction of incoherent probabilities (Hope, Korb & Nicholson, 2002)
 - Minimizing squared deviations from original estimates
- Automated maxentropy fill of CPTs (Hope, Korb & Nicholson, 2002)
- Automated normalization of CPTs (Hope, Korb & Nicholson, 2002)
- Use of lotteries to force estimates (also useful for utility elicitation)

Local model structure

Not every cell in CPT is independent from every other cell. Examples:

- Deterministic nodes
 - It is possible to have nodes where the value of a child is exactly specified (logically or numerically) by its parents

- Linear relationships:

$$X_i = a_0 X_0 + \dots a_n X_n + \epsilon_i$$

- Logit model (binary, 2 parents):

$$\log \frac{P(X_2|X_0, X_1)}{P(\neg X_2|X_0, X_1)} = a + bX_0 + cX_1 + dX_1X_2$$

- Partitions of parent state space
- Independence of causal influence
- Contingent substructures

Decision Analysis

Since 1970s there have been nice software packages for decision analysis:

- Eliciting actions
- Eliciting utilities
- Eliciting probabilities
- Building decision trees
- Sensitivity analysis, etc.

See: Raiffa's *Intro to Decision Analysis* (an excellent book!)

Main differences from KEBN:

- *Scale*: tens vs thousands of parms!!
- *Structure*: trees reflect state-action combinations, not causal structure, prediction, intervention

Eliciting Decision Networks

- Action nodes: What actions can be taken in domain?
- Utility node(s):
 - What unit(s) will “utile” be measured in?
 - Are there difference aspects to the utility that should each be represented in a separate utility node?
- Graph structure:
 - Which variables can decision/actions affects?
 - Does the action/decision affect the utility?
 - What are the outcome variables that there are preferences about?

Model Evaluation

- Elicitation review
 - Review variable and value definition
 - * clarity test, agreement on definitions, consistency
 - Review graph and local model structure
 - Review probabilities
 - * compare probabilities with each other
- Sensitivity analysis (Laskey, 1993)
 - Measures effect of one variable on another
- Case-based evaluation
 - Run model on test of test cases
 - Compare with expert judgement or “ground truth”
- Validation methods using data (if available)
 - Predictive Accuracy
 - Expected value
 - Kullback-Leibler divergence
 - (Bayesian) Information reward

The need to prototype!

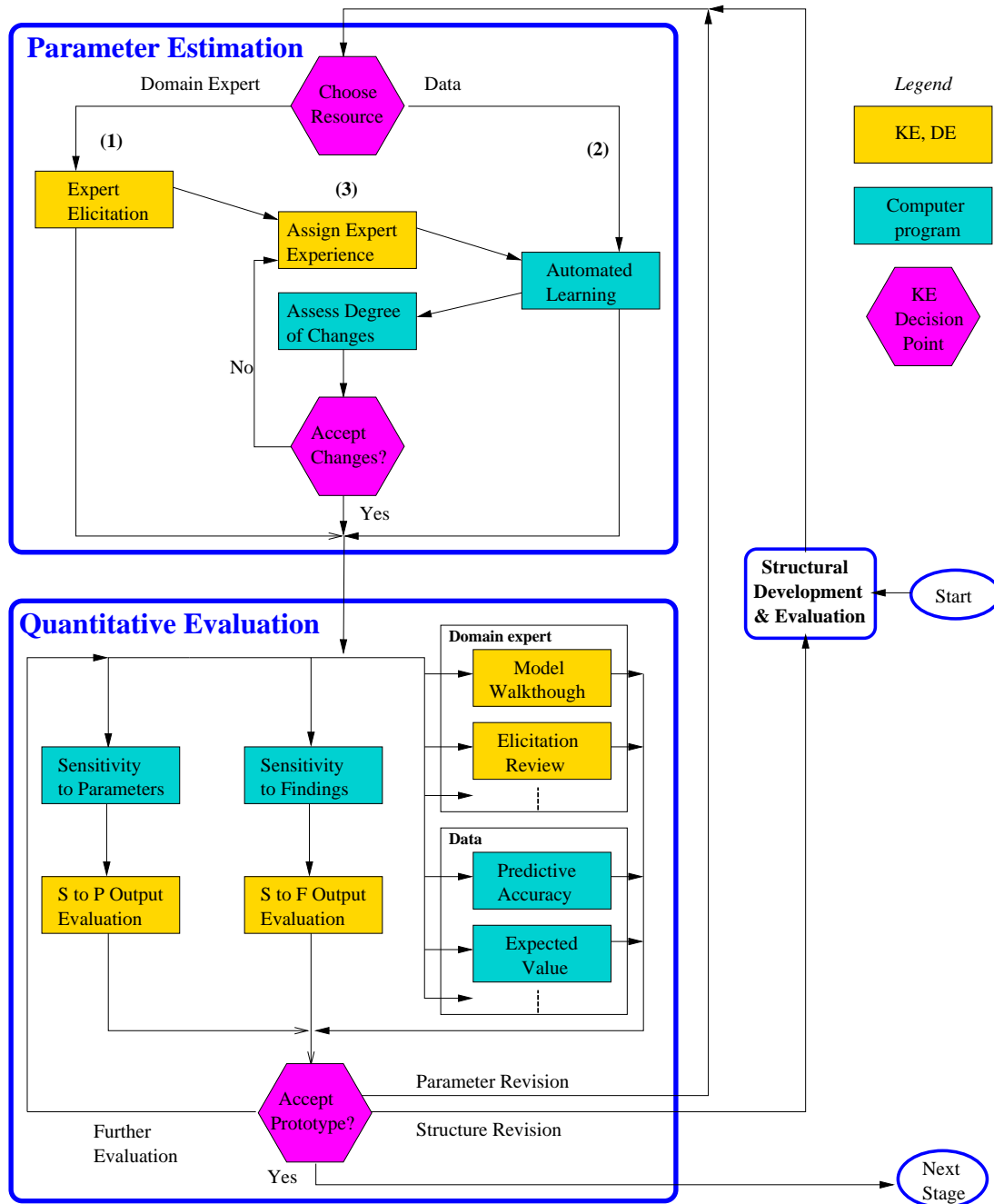
Why prototype?

- It's just the best software development process overall (Brooks). Organic growth of software:
 - tracks the specs
 - has manageable size (at least initially)
- Attacks the comprehensiveness vs. intelligibility trade-off from the right starting point.
- Few off-the-shelf models; prototyping helps us fill in the gaps, helps write the specs

Prototypes

- Initial prototypes minimize risk
 - Don't oversell result
 - Employ available capabilities
 - Simplify variables, structure, questions answered
 - Provide working product for assessment
- Incremental prototypes
 - Simple, quick extension to last
 - Attacks high priority subset of difficult issues
 - Helps refine understanding of requirements/approach

KEBN methodologies



KEBN Summary

- Various BN structures are available to compactly and accurately represent certain types of domain features.
- There is an interplay between elements of the KE process: variable choice, graph structure and parameters.
- No standard knowledge engineering process exists as yet.
- Integration of expert elicitation and automated methods still in early stages.
- There are few existing tools for supporting the BN KE process.
 - We at Monash are developing some!

Monash BN Applications: Overview

- User modelling (plan recognition in a MUD, web page pre-fetching): Zukerman, Albrecht, Nicholson (1997-2001)
- Ambulation monitoring and fall detection: Nicholson, Brown (Monash Biomedical Engineering), Honours projects 1997, 2000
- Seabreeze prediction: Nicholson, Korb, Bureau of Meteorology, 2001 Honours projects
- Intelligent tutoring for decimal understanding: Nicholson, Boneh, University of Melbourne (1999-2003)
- NAG (Nice Argument Generator): Zukerman, Korb
- Bayesian Poker: Korb, Nicholson, Honours projects 1993,1994,1995,2001,2003
- SARBayes: Twardy, Korb, Albrecht, Victorian Search and Rescue, 2001 Honours project

Monash BN Applications (cont.)

- Ecological risk assessment:
 - Nicholson, Korb, Pollino (Monash Centre for Water Studies), 2003-2005 Native Fish abundance in Goulburn Water
 - Predicting recreational water quality: Twardy, Nicholson, NSW EPA, 2003 Honours project
 - Tropical seagrass in great barrier reef: Nicholson, Thomas (Monash Centre for Water Studies), 2004-2006
- Change impact analysis in software architecture design: Nicholson, Tang, Jin, Han (Swinburne)

References

Introduction to Bayesian AI

- T. Bayes (1764) "An Essay Towards Solving a Problem in the Doctrine of Chances." *Phil Trans of the Royal Soc of London*. Reprinted in *Biometrika*, 45 (1958), 296-315.
- B. Buchanan and E. Shortliffe (eds.) (1984) *Rule-Based Expert Systems: The MYCIN Experiments of the Stanford Heuristic Programming Project*. Addison-Wesley.
- R. Cox (1946). Probability, frequency and reasonable expectation. *American Journal of Physics*, 14, 1-13, 1946.
- B. de Finetti (1964) "Foresight: Its Logical Laws, Its Subjective Sources," in Kyburg and Smokler (eds.) *Studies in Subjective Probability*. NY: Wiley.
- G. Gigerenzer and U. Hoffrage (1995). "How to improve Bayesian reasoning without instruction: Frequency formats." *Psychological Review*, 102, 684-704.
- A. Hajek (2002). "Scotching Dutch Books." Talk at Monash University, July, 2002.
- D. Heckerman (1986) "Probabilistic Interpretations for MYCIN's Certainty Factors," in L.N. Kanal and J.F. Lemmer (eds.) *Uncertainty in Artificial Intelligence*. North-Holland.

- E. J. Horvitz, J. S. Breese, and M. Henrion (1988). Decision theory in expert systems and artificial intelligence. *International Journal of Approximate Reasoning*, 2, 247-302.
- C. Howson and P. Urbach (1993) *Scientific Reasoning: The Bayesian Approach*. Open Court.
A MODERN REVIEW OF BAYESIAN THEORY.
- D. Hume (1737) *A Treatise of Human Nature*.
- D. Kahneman, P. Slovic and A. Tversky (eds.) (1982) *Judgment under uncertainty: Heuristics and biases*. Cambridge.
- K.B. Korb (1995) "Inductive learning and defeasible inference," *Jrn for Experimental and Theoretical AI*, 7, 291-324.
- R. Neapolitan (1990) *Probabilistic Reasoning in Expert Systems*. Wiley.
CHAPTERS 1, 2 AND 4 COVER SOME OF THE RELEVANT HISTORY.
- J. Pearl (1988) *Probabilistic Reasoning in Intelligent Systems*, Morgan Kaufmann.
- F. P. Ramsey (1926/1931) "Truth and Probability" in *The Foundations of Mathematics and Other Essays*. NY: Humanities Press.
THE ORIGIN OF MODERN BAYESIANISM. INCLUDES LOTTERY-BASED ELICITATION AND DUTCH-BOOK ARGUMENTS FOR THE USE OF PROBABILITIES.

R. Reiter (1980) "A logic for default reasoning," *Artificial Intelligence*, 13, 81-132.

J. von Neumann and O. Morgenstern (1947) *Theory of Games and Economic Behavior*, 2nd ed. Princeton Univ.

STANDARD REFERENCE ON ELICITING UTILITIES VIA LOTTERIES.

Bayesian Networks

E. Charniak (1991) "Bayesian Networks Without Tears", *Artificial Intelligence Magazine*, pp. 50-63, Vol 12.

AN ELEMENTARY INTRODUCTION.

G.F. Cooper (1990) The computational complexity of probabilistic inference using belief networks. *Artificial Intelligence*, 42, 393-405.

R. G. Cowell, A. Philip Dawid, S. L. Lauritzen and D. J. Spiegelhalter (1999) *Probabilistic networks and expert systems*. New York: Springer.

TECHNICAL SURVEY OF BAYESIAN NET TECHNOLOGY, INCLUDING LEARNING BAYESIAN NETS.

D. D'Ambrosio (1999) "Inference in Bayesian Networks". *Artificial Intelligence Magazine*, Vol 20, No. 2.

A. P. Dawid (1998) Conditional independence. In *Encyclopedia of Statistical Sciences, Update Volume 2*. New York: Wiley Interscience.

P. Haddaway (1999) "An Overview of Some Recent Developments in Bayesian Problem-Solving Techniques". *Artificial Intelligence Magazine*, Vol 20, No. 2.

R.A. Howard & J.E. Matheson (1981) Influence Diagrams.
In Howard and Matheson (eds.) *Readings in the Principles and Applications of Decision Analysis*. Menlo Park, Calif: Strategic Decisions Group.

F. V. Jensen (1996) *An Introduction to Bayesian Networks*, Springer.

K.B. Korb, L.R. Hope, A.E. Nicholson and K. Axnick.
Varieties of Causal Intervention In C. Zhang, H.W. Guesgen and W.K. Yeap (eds), *Lecture Notes in Artificial Intelligence, (Proceedings of the 8th Pacific Rim International Conference on Artificial Intelligence, [PRICAI 2004], Auckland, New Zealand, August 2004)*, Springer-Verlag, Berlin, Germany, Vol 3157, pp 322-331.

R. Neapolitan (1990) *Probabilistic Reasoning in Expert Systems*. Wiley.

SIMILAR COVERAGE TO THAT OF PEARL; MORE EMPHASIS ON PRACTICAL ALGORITHMS FOR NETWORK UPDATING.

J. Pearl (1988) *Probabilistic Reasoning in Intelligent Systems*, Morgan Kaufmann.

THIS IS THE CLASSIC TEXT INTRODUCING BAYESIAN NETWORKS TO THE AI COMMUNITY.

J. Pearl (2000) *Causality*. Cambridge University.

Poole, D., Mackworth, A., and Goebel, R. (1998)
Computational Intelligence: a logical approach. Oxford University Press.

Russell & Norvig (1995) *Artificial Intelligence: A Modern Approach*, Prentice Hall.

J. Whittaker (1990) *Graphical models in applied multivariate statistics*. Wiley.

Learning Bayesian Networks

H. Blalock (1964) *Causal Inference in Nonexperimental Research*. University of North Carolina.

R. Bouckear (1994) *Probabilistic network construction using the minimum description length principle*. Technical Report RUU-CS-94-27, Dept of Computer Science, Utrecht University.

C. Bouillier, N. Friedman, M. Goldszmidt, D. Koller (1996) "Context-specific independence in Bayesian networks," in Horvitz & Jensen (eds.) *UAI 1996*, 115-123.

G. Brightwell and P. Winkler (1990) *Counting linear extensions is #P-complete*. Technical Report DIMACS 90-49, Dept of Computer Science, Rutgers Univ.

W. Buntine (1991) "Theory refinement on Bayesian networks," in D'Ambrosio, Smets and Bonissone (eds.) *UAI 1991*, 52-69.

W. Buntine (1996) "A Guide to the Literature on Learning Probabilistic Networks from Data," *IEEE Transactions on Knowledge and Data Engineering*, 8, 195-210.

D.M. Chickering (1995) "A Transformational Characterization of Equivalent Bayesian Network Structures," in P. Besnard and S. Hanks (eds.) *Proceedings of the Eleventh Conference on Uncertainty in Artificial Intelligence* (pp. 87-98). San Francisco: Morgan Kaufmann.

STATISTICAL EQUIVALENCE.

Chickering, D. M. (2002). "Optimal structure identification with greedy search." *Journal of Machine Learning Research*, 3, 507-559.

G.F. Cooper and E. Herskovits (1991) "A Bayesian Method for Constructing Bayesian Belief Networks from Databases," in D'Ambrosio, Smets and Bonissone (eds.) *UAI 1991*, 86-94.

G.F. Cooper and E. Herskovits (1992) "A Bayesian Method for the Induction of Probabilistic Networks from Data," *Machine Learning*, 9, 309-347.

AN EARLY BAYESIAN CAUSAL DISCOVERY METHOD.

H. Dai, K.B. Korb, C.S. Wallace and X. Wu (1997) "A study of casual discovery with weak links and small samples." *Proceedings of the Fifteenth International Joint Conference on Artificial Intelligence (IJCAI)*, pp. 1304-1309. Morgan Kaufmann.

N. Friedman (1997) "The Bayesian Structural EM Algorithm," in D. Geiger and P.P. Shenoy (eds.) *Proceedings of the Thirteenth Conference on Uncertainty in Artificial Intelligence* (pp. 129-138). San Francisco: Morgan Kaufmann.

D. Geiger and D. Heckerman (1994) "Learning Gaussian networks," in Lopes de Mantras and Poole (eds.) *UAI 1994*, 235-243.

D. Heckerman and D. Geiger (1995) "Learning Bayesian networks: A unification for discrete and Gaussian domains," in Besnard and Hanks (eds.) *UAI 1995*, 274-284.

D. Heckerman, D. Geiger, and D.M. Chickering (1995) "Learning Bayesian Networks: The Combination of Knowledge and Statistical Data," *Machine Learning*, 20, 197-243.

BAYESIAN LEARNING OF STATISTICAL EQUIVALENCE CLASSES.

K. Korb (1999) "Probabilistic Causal Structure" in H. Sankey (ed.) *Causation and Laws of Nature: Australasian Studies in History and Philosophy of Science 14*. Kluwer Academic.

INTRODUCTION TO THE RELEVANT PHILOSOPHY OF CAUSATION FOR LEARNING BAYESIAN NETWORKS.

P. Krause (1998) *Learning Probabilistic Networks*.

[http : //www.auai.org/bayes_USKrause.ps.gz](http://www.auai.org/bayes_USKrause.ps.gz)

BASIC INTRODUCTION TO BNS, PARAMETERIZATION AND LEARNING CAUSAL STRUCTURE.

W. Lam and F. Bacchus (1993) "Learning Bayesian belief networks: An approach based on the MDL principle," *Jrn Comp Intelligence*, 10, 269-293.

D. Madigan, S.A. Andersson, M.D. Perlman & C.T. Volinsky (1996) "Bayesian model averaging and model selection for Markov equivalence classes of acyclic digraphs," *Comm in Statistics: Theory and Methods*, 25, 2493-2519.

D. Madigan and A. E. Raftery (1994) "Model selection and accounting for model uncertainty in graphical models using Occam's window," *Jrn AMer Stat Assoc*, 89, 1535-1546.

Meek, C. (1996). *Graphical models: Selectiong causal and statistical models*. PhD disseration, Philosophy, Carnegie Mellon University.

N. Metropolis, A.W. Rosenbluth, M.N. Rosenbluth, A.H. Teller and E. Teller (1953) "Equations of state calculations by fast computing machines," *Jrn Chemical Physics*, 21, 1087-1091.

J.R. Neil and K.B. Korb (1999) "The Evolution of Causal Models: A Comparison of Bayesian Metrics and Structure Priors," in N. Zhong and L. Zhous (eds.) *Methodologies for Knowledge Discovery and Data Mining: Third Pacific-Asia Conference* (pp. 432-437). Springer Verlag.

GENETIC ALGORITHMS FOR CAUSAL DISCOVERY;
STRUCTURE PRIORS.

J.R. Neil, C.S. Wallace and K.B. Korb (1999) "Learning Bayesian networks with restricted causal interactions," in Laskey and Prade (eds.) *UAI 99*, 486-493.

- J. Rissanen (1978) "Modeling by shortest data description," *Automatica*, 14, 465-471.
- R.W. Robinson (1977). "Counting unlabelled acyclic diagraphs," in C.H.C.Little (Ed.), *Lecture notes in mathematics 62: Combinatorial mathematics V*. New York: Springer-Verlag.
- H. Simon (1954) "Spurious Correlation: A Causal Interpretation," *Jrn Amer Stat Assoc*, 49, 467-479.
- D. Spiegelhalter & S. Lauritzen (1990) "Sequential Updating of Conditional Probabilities on Directed Graphical Structures," *Networks*, 20, 579-605.
- P. Spirtes, C. Glymour and R. Scheines (1990) "Causality from Probability," in J.E. Tiles, G.T. McKee and G.C. Dean *Evolving Knowledge in Natural Science and Artificial Intelligence*. London: Pitman. AN ELEMENTARY INTRODUCTION TO STRUCTURE LEARNING VIA CONDITIONAL INDEPENDENCE.
- P. Spirtes, C. Glymour and R. Scheines (1993) *Causation, Prediction and Search: Lecture Notes in Statistics 81*. Springer Verlag.
A THOROUGH PRESENTATION OF THE ORTHODOX STATISTICAL APPROACH TO LEARNING CAUSAL STRUCTURE.
- J. Suzuki (1996) "Learning Bayesian Belief Networks Based on the Minimum Description Length Principle," in L. Saitta (ed.) *Proceedings of the Thirteenth International Conference on Machine Learning* (pp. 462-470). San Francisco: Morgan Kaufmann.

T.S. Verma and J. Pearl (1991) “Equivalence and Synthesis of Causal Models,” in P. Bonissone, M. Henrion, L. Kanal and J.F. Lemmer (eds) *Uncertainty in Artificial Intelligence 6* (pp. 255-268). Elsevier.

THE GRAPHICAL CRITERION FOR STATISTICAL EQUIVALENCE.

C.S. Wallace and D. Boulton (1968) “An information measure for classification,” *Computer Jrn*, 11, 185-194.

C.S. Wallace and P.R. Freeman (1987) “Estimation and inference by compact coding,” *Jrn Royal Stat Soc (Series B)*, 49, 240-252.

C. S. Wallace and K. B. Korb (1999) “Learning Linear Causal Models by MML Sampling,” in A. Gammernan (ed.) *Causal Models and Intelligent Data Management*. Springer Verlag.

SAMPLING APPROACH TO LEARNING CAUSAL MODELS; DISCUSSION OF STRUCTURE PRIORS.

C. S. Wallace, K. B. Korb, and H. Dai (1996) “Causal Discovery via MML,” in L. Saitta (ed.) *Proceedings of the Thirteenth International Conference on Machine Learning* (pp. 516-524). San Francisco: Morgan Kaufmann.

INTRODUCES AN MML METRIC FOR CAUSAL MODELS.

S. Wright (1921) “Correlation and Causation,” *Jrn Agricultural Research*, 20, 557-585.

S. Wright (1934) “The Method of Path Coefficients,” *Annals of Mathematical Statistics*, 5, 161-215.

BN Knowledge Engineering

Boneh, T. (2002) “Visualisation of structural dependencies for Bayesian Network Knowledge Engineering”, Masters thesis, University of Melbourne.

C. Boutilier, N. Friedman, M. Goldszmidt, D. Koller (1996) “Context-specific independence in Bayesian networks,” in Horvitz & Jensen (eds.) *UAI 1996*, 115-123.

Druzdzel, M.J. and van der Gaag, L.C. (1995), “Elicitation of probabilities for belief networks: Combining qualitative and quantitative information”, Besnard & Hanks (eds) *UAI95*, pp. 141-148.

M.J. Druzdzel and L.C. van der Gaag (Eds.) (2000) “Building probabilistic networks: Where do the numbers come from?” Editors Introduction to Special Section *IEEE Trans. on Knowledge and Data Engineering*, 12(4), pp 481-486.

D. Heckerman (1990), “Probabilistic Similarity Networks”, *Networks*, 20, pp. 607-636.

Laskey, K.B. (1993) “Sensitivity Analysis for Probability Assessments in Bayesian Networks”, in Heckerman & Mamdani (eds.) *UAI93*, pp. 136-142.

Laskey, K.B. and Mahoney, S.M. (1997) “Network Fragments: Representing Knowledge for Constructing Probabilistic Models”, in Geiger & Shenoy (eds.) *UAI97*, pp. 334-341.

Laskey, K.B. (1999) A Full-day tutorial given at *UAI99*.

<http://ite.gmu.edu/klaskey/papers/uai99tutorial.pdf>

Laskey, K.B. and Mahoney, S.M. (2000) "Network Engineering for Agile Belief Network Models", *IEEE: Transactions on Knowledge and Data Engineering*, 12(4), pp. 487-498.

Mahoney, S.M. and Laskey, K.B. (1996) "Network Engineering for Complex Belief Networks", in Horvitz & Jensen (eds.) *UAI96*, pp. 389-396.

Monti, S. and Carenini, G. (2000) "Dealing with the Expert Inconsistency in Probability Elicitation", *IEEE: Transactions on Knowledge and Data Engineering*, 12(4), pp. 499-508.

Nikovski, D. (2000) "Constructing Bayesian Networks for Medical Diagnosis from Incomplete and Partially Correct Statistics", *IEEE: Transactions on Knowledge and Data Engineering*, 12(4), pp. 509-516.

M. Pradham, G. Provan, B. Middleton and M. Henrion (1994) "Knowledge engineering for large belief networks", In de Mantras & Poole (eds.) *UAI94*, pp. 484-490.

S. Srinivas (1993) "A Generalization of the Noisy-OR model", In Heckerman & Mamdani (eds.) *UAI93* pp. 208-215.

L.C. van der Gaag, S. Renooij, C.L.M. Witteman, B.M.P. Aleman, B.G. Taal (1999) "How to Elicit Many Probabilities", Laskey & Prade (eds) *UAI99*, pp. 647-654.

O. Woodberry, A.E. Nicholson, K.B. Korb, and C. Pollino (2004). Parameterising Bayesian Networks. in G. I. Webb and X. Yu (eds). *Lecture Notes in Artificial Intelligence, (Proceedings of the 17th Australian Joint Conference on Advances in Artificial Intelligence [AI'04], Cairns, Australia, 4-6 December 2004)*, Springer-Verlag, Berlin, Germany, Vol 3339, pp 1101-1107.