**Lecture 5**
# Baysian Networks

Marco Chiarandini

Department of Mathematics & Computer Science
University of Southern Denmark

Slides by Stuart Russell and Peter Norvig

# Course Overview

- ✔ Introduction
    - ✔ Artificial Intelligence
    - ✔ Intelligent Agents
- ✔ Search
    - ✔ Uninformed Search
    - ✔ Heuristic Search
- Uncertain knowledge and Reasoning
    - Probability and Bayesian approach
    - Bayesian Networks
    - Hidden Markov Chains
    - Kalman Filters

- Learning
    - Supervised
      Learning Bayesian Networks, Neural Networks
    - Unsupervised
      EM Algorithm
- Reinforcement Learning
- Games and Adversarial Search
    - Minimax search and Alpha-beta pruning
    - Multiagent search
- Knowledge representation and Reasoning
    - Propositional logic
    - First order logic
    - Inference
    - Plannning

# Summary

Probability is a rigorous formalism for uncertain knowledge

Joint probability distribution specifies probability of every atomic event
Queries can be answered by summing over atomic events

For nontrivial domains, we must find a way to reduce the joint size

Independence and conditional independence provide the tools

# Outline

# Outline

◇ Syntax
◇ Semantics
◇ Parameterized distributions

# Bayesian networks

### Definition

A simple, graphical notation for conditional independence assertions and hence for compact specification of full joint distributions

Syntax:
- a set of nodes, one per variable
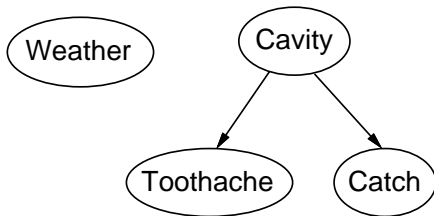- a directed, acyclic graph (link $\approx$ "directly influences")
- a conditional distribution for each node given its parents:
  $\Pr(X_i \mid Parents(X_i))$

In the simplest case, conditional distribution represented as a conditional probability table (CPT) giving the distribution over $X_i$ for each combination of parent values

# Example

Topology of network encodes conditional independence assertions:



*Weather* is independent of the other variables
*Toothache* and *Catch* are conditionally independent given *Cavity*
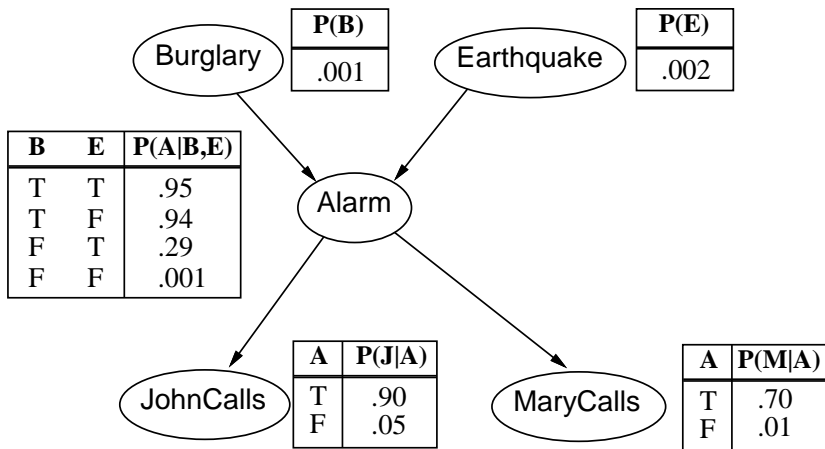
# Example

I'm at work, neighbor John calls to say my alarm is ringing, but neighbor Mary doesn't call. Sometimes it's set off by minor earthquakes. Is there a burglar?

Variables: *Burglar*, *Earthquake*, *Alarm*, *JohnCalls*, *MaryCalls*

Network topology reflects "causal" knowledge:
  – A burglar can set the alarm off
  – An earthquake can set the alarm off
  – The alarm can cause Mary to call
  – The alarm can cause John to call

**P(B)**

| |
|---|
| .001 |

**P(E)**

| |
|---|
| .002 |

| B | E | P(A\|B,E) |
|---|---|---|
| T | T | .95 |
| T | F | .94 |
| F | T | .29 |
| F | F | .001 |

| A | P(J\|A) |
|---|---|
| T | .90 |
| F | .05 |

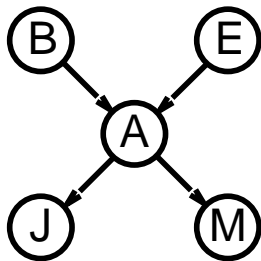| A | P(M\|A) |
|---|---|
| T | .70 |
| F | .01 |

# Compactness

A CPT for Boolean $X_i$ with $k$ Boolean parents has $2^k$ rows for the combinations of parent values

Each row requires one number $p$ for $X_i = true$
(the number for $X_i = false$ is just $1 - p$)
If each variable has no more than $k$ parents,
the complete network requires $O(n \cdot 2^k)$ numbers

I.e., grows linearly with $n$, vs. $O(2^n)$ for the full joint distribution

For burglary net, $1 + 1 + 4 + 2 + 2 = 10$ numbers
(vs. $2^5 - 1 = 31$)

# Global semantics

"Global" semantics defines the full joint distribution as the product of the local conditional distributions:
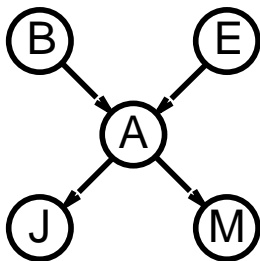
$$P(x_1, \ldots, x_n) = \prod_{i=1}^{n} P(x_i \mid parents(X_i))$$

e.g., $P(j \wedge m \wedge a \wedge \neg b \wedge \neg e)$

$$= \quad P(j \mid a)P(m \mid a)P(a \mid \neg b, \neg e)P(\neg b)P(\neg e)$$

$$= \quad 0.9 \times 0.7 \times 0.001 \times 0.999 \times 0.998$$

$$\approx \quad 0.00063$$

# Constructing Bayesian networks

Need a method such that a series of locally testable assertions of conditional independence guarantees the required global semantics

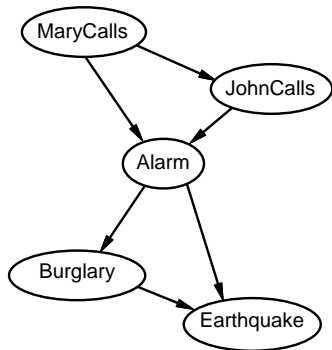- Choose an ordering of variables $X_1, \ldots, X_n$
- For $i = 1$ to $n$
  add $X_i$ to the network
  select parents from $X_1, \ldots, X_{i-1}$ such that
  $\Pr(X_i \mid Parents(X_i)) = \Pr(X_i \mid X_1, \ldots, X_{i-1})$

This choice of parents guarantees the global semantics:

$$
\begin{aligned}
\Pr(X_1, \ldots, X_n) &= \prod_{i=1}^{n} \Pr(X_i \mid X_1, \ldots, X_{i-1}) \quad \text{(chain rule)} \\
&= \prod_{i=1}^{n} \Pr(X_i \mid Parents(X_i)) \quad \text{(by construction)}
\end{aligned}
$$

# Example

Suppose we choose the ordering $M$, $J$, $A$, $B$, $E$



$P(J \mid M) = P(J)$?   No
$P(A \mid J, M) = P(A \mid J)$?
$P(A \mid J, M) = P(A)$?   No
$P(B \mid A, J, M) = P(B \mid A)$?   Yes
$P(B \mid A, J, M) = P(B)$?   No
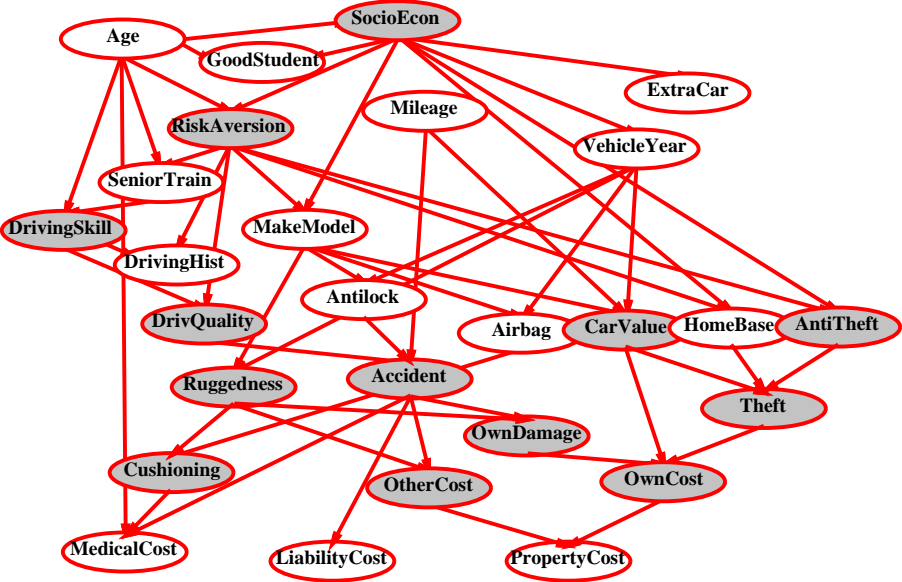$P(E \mid B, A, J, M) = P(E \mid A)$?   No
$P(E \mid B, A, J, M) = P(E \mid A, B)$?
Yes
Deciding conditional independence is
hard in noncausal directions
(Causal models and conditional
independence seem hardwired for
humans!)
Assessing conditional probabilities is
hard in noncausal directions
Network is less compact:
$1 + 2 + 4 + 2 + 4 = 13$ numbers needed

14

# Example: Car insurance

# Compact conditional distributions

CPT grows exponentially with number of parents
CPT becomes infinite with continuous-valued parent or child

Solution:
canonical distributions that are defined compactly

Deterministic nodes are the simplest case:
$X = f(Parents(X))$ for some function $f$

E.g., Boolean functions
$NorthAmerican \Leftrightarrow Canadian \lor US \lor Mexican$

E.g., numerical relationships among continuous variables

$$\frac{\partial Level}{\partial t} = \text{inflow + precipitation - outflow - evaporation}$$

# Compact conditional distributions contd.

Noisy-OR distributions model multiple noninteracting causes
  1) Parents $U_1 \ldots U_k$ include all causes (can add leak node)
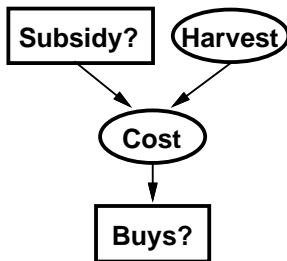  2) Independent failure probability $q_i$ for each cause alone

$$\implies P(X \mid U_1 \ldots U_j, \neg U_{j+1} \ldots \neg U_k) = 1 - \prod_{i=1}^{j} q_i$$

| Cold | Flu | Malaria | P(Fever) | P(¬Fever) |
|------|-----|---------|----------|-----------|
| F | F | F | 0.0 | 1.0 |
| F | F | T | 0.9 | 0.1 |
| F | T | F | 0.8 | 0.2 |
| F | T | T | 0.98 | $0.02 = 0.2 \times 0.1$ |
| T | F | F | 0.4 | 0.6 |
| T | F | T | 0.94 | $0.06 = 0.6 \times 0.1$ |
| T | T | F | 0.88 | $0.12 = 0.6 \times 0.2$ |
| T | T | T | 0.988 | $0.012 = 0.6 \times 0.2 \times 0.1$ |

Number of parameters **linear** in number of parents

# Hybrid (discrete+continuous) networks

Discrete (*Subsidy?* and *Buys?*); continuous (*Harvest* and *Cost*)



Option 1: discretization—possibly large errors, large CPTs
Option 2: finitely parameterized canonical families
1) Continuous variable, discrete+continuous parents (e.g., *Cost*)
2) Discrete variable, continuous parents (e.g., *Buys?*)

# Continuous child variables

Need one conditional density function for child variable given continuous parents, for each possible assignment to discrete parents

Most common is the linear Gaussian model, e.g.,:

$$P(Cost = c \mid Harvest = h, Subsidy = true)$$
$$= N(a_t h + b_t, \sigma_t)$$
$$= \frac{1}{\sigma_t \sqrt{2\pi}} exp\left(-\frac{1}{2}\left(\frac{c - (a_t h + b_t)}{\sigma_t}\right)^2\right)$$

Mean *Cost* varies linearly with *Harvest*, variance is fixed

⤳Linear variation is unreasonable over the full range
   but works OK if the **likely** range of *Harvest* is narrow
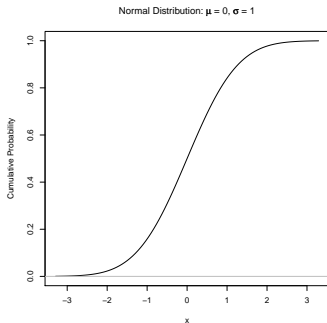
# Continuous child variables

P(Cost|Harvest,Subsidy?=true)

All-continuous network with linear Gaussian distributions
$\implies$ full joint distribution is a multivariate Gaussian

Discrete+continuous linear Gaussian network is a conditional Gaussian
network i.e., a multivariate Gaussian over all continuous variables for each
combination of discrete variable values

# Discrete variable w/ continuous parents

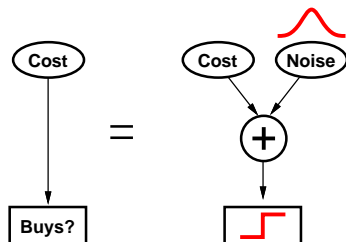Probability of *Buys?* given *Cost* should be a "soft" threshold:



Normal Distribution: $\mu = 0$, $\sigma = 1$

Probit distribution uses integral of Gaussian:
$$\Phi(x) = \int_{-\infty}^{x} N(0,1)(x)dx$$
$$P(Buys? = true \mid Cost = c) = \Phi((-c + \mu)/\sigma)$$

# Why the probit?

1. It's sort of the right shape
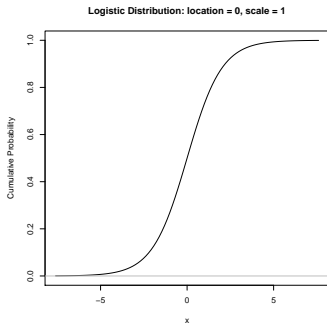2. Can be viewed as hard threshold whose location is subject to noise

# Discrete variable contd.

Sigmoid (or logit) distribution also used in neural networks:

$$P(Buys? = true \mid Cost = c) = \frac{1}{1 + exp(-2\frac{-c+\mu}{\sigma})}$$

Sigmoid has similar shape to probit but much longer tails:



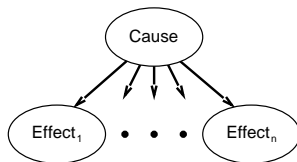Logistic Distribution: location = 0, scale = 1
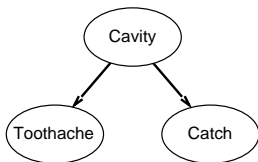
# Summary

- Bayes nets provide a natural representation for (causally induced) conditional independence

- Topology + CPTs = compact representation of joint distribution

- Generally easy for (non)experts to construct

- Canonical distributions (e.g., noisy-OR) = compact representation of CPTs

- Continuous variables $\implies$ parameterized distributions (e.g., linear Gaussian)

# Bayes' Rule and conditional independence

$\Pr(Cavity \mid toothache \land catch)$

$= \alpha \Pr(toothache \land catch \mid Cavity) \Pr(Cavity)$

$= \alpha \Pr(toothache \mid Cavity) \Pr(catch \mid Cavity) \Pr(Cavity)$
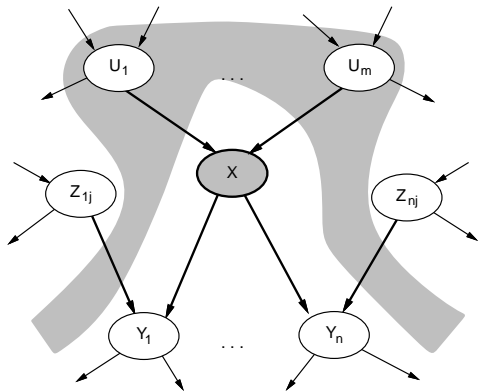
This is an example of a naive Bayes model:

$$\Pr(Cause, Effect_1, \ldots, Effect_n) = \Pr(Cause) \prod_i \Pr(Effect_i \mid Cause)$$



Total number of parameters is **linear** in $n$

# Local semantics

Local semantics: each node is conditionally independent
of its nondescendants given its parents



Theorem: Local semantics $\Leftrightarrow$ global semantics

# Markov blanket

Each node is conditionally independent of all others given its
Markov blanket: parents + children + children's parents