

PROBABILISTIC REASONING

Outline

- Representing Knowledge in an uncertain domain
- Semantics of Bayesian networks
- Efficient representation of conditional distributions
- Exact/Approximate inference in Bayesian networks
- Other approaches to uncertain reasoning

Based on the textbook by Stuart Russell & Peter Norvig:

Artificial Intelligence, A Modern Approach (2nd Edition)

Chapter 14; excluding Section 14.6

Bayesian Networks: Syntax

Definition. A belief network is a *directed acyclic graph* (DAG) $G = \langle \{X_1, \dots, X_n\}, E \rangle$ where

1. nodes X_1, \dots, X_n are discrete/continuous random variables,
2. the set of *arrows* (or links)

$$E \subseteq \{X_1, \dots, X_n\}^2 = \{ \langle X_i, X_j \rangle \mid 1 \leq i \leq n \text{ and } 1 \leq j \leq n \},$$
3. an arrow $\langle X, Y \rangle \in E$ of G represents a *direct influence* relationship between the variables X and Y , and
4. each node X is assigned a completely specified probability distribution $\mathbf{P}(X|\text{Parents}(X))$ where

$$\text{Parents}(X) = \{Y \mid \langle Y, X \rangle \in E\}.$$

1. REPRESENTING KNOWLEDGE IN AN UNCERTAIN DOMAIN

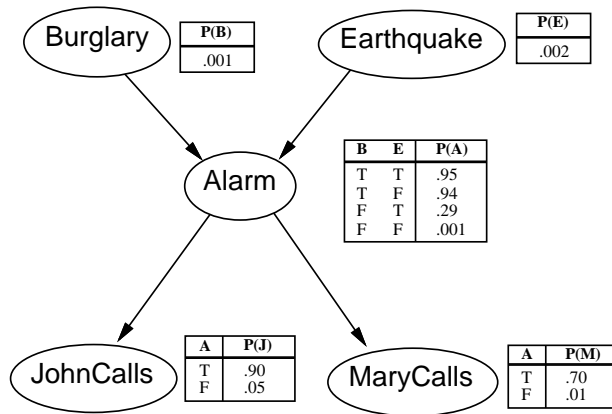
- Conditional independence relations provide means to simplify probabilistic representations of the world.
- A **Bayesian network** is a data structure representing the dependencies among variables X_1, \dots, X_n of a given domain.
- As a result, a compact specification of the full joint probability distribution $\mathbf{P}(X_1, \dots, X_n)$ is obtained.
- Bayesian networks are also called *belief networks*, *probabilistic networks*, *causal networks* or *knowledge maps*.

Example. Consider a network based on five Boolean random variables:

1. *Burglary* = "a burglar enters our home".
2. *Earthquake* = "an earthquake occurs".
3. *Alarm* = "our burglar alarm goes off".
The alarm is fairly reliable at detecting a burglary, but may occasionally respond to minor earthquakes.
4. *JohnCalls* = "Our neighbor John calls and reports an alarm."
He always calls when he hears the alarm, but sometimes confuses telephone ringing with the alarm.
5. *MaryCalls* = "Our neighbor Mary calls and reports an alarm". She likes loud music and sometimes misses the alarm altogether.

Shorthands B , E , A , J , and M are also introduced for these variables.

- The relationships of the variables are given as a Bayesian network.
- The probability distributions $\mathbf{P}(X \mid \text{Parents}(X))$ associated with variables X are given as *conditional probability tables* (CPTs).



© 2006 TKK / Laboratory for Theoretical Computer Science

Conditional Independence Revisited

Definition. Let $P(\psi) > 0$. Sentences ϕ_1 and ϕ_2 are *conditionally independent given ψ* $\iff P(\phi_1 \wedge \phi_2 \mid \psi) = P(\phi_1 \mid \psi)P(\phi_2 \mid \psi)$.

Proposition. If $P(\psi) > 0$, $P(\phi_1 \wedge \psi) > 0$, and $P(\phi_2 \wedge \psi) > 0$, then ϕ_1 and ϕ_2 are conditionally independent given $\psi \iff P(\phi_1 \mid \phi_2 \wedge \psi) = P(\phi_1 \mid \psi)$ and $P(\phi_2 \mid \phi_1 \wedge \psi) = P(\phi_2 \mid \psi)$ hold.

Proof. For the former equation, we note that

$$\begin{aligned} P(\phi_1 \wedge \phi_2 \mid \psi) &= P(\phi_1 \mid \psi)P(\phi_2 \mid \psi) \\ \iff \frac{P(\phi_1 \wedge \phi_2 \wedge \psi)}{P(\psi)} &= \frac{P(\phi_1 \wedge \psi)}{P(\psi)} \cdot \frac{P(\phi_2 \wedge \psi)}{P(\psi)} \\ \iff P(\phi_1 \wedge \phi_2 \wedge \psi)P(\psi) &= P(\phi_1 \wedge \psi)P(\phi_2 \wedge \psi) \\ \iff P(\phi_1 \mid \phi_2 \wedge \psi) &= \frac{P(\phi_1 \wedge \phi_2 \wedge \psi)}{P(\phi_2 \wedge \psi)} = \frac{P(\phi_1 \wedge \psi)}{P(\psi)} = P(\phi_1 \mid \psi). \end{aligned}$$

© 2006 TKK / Laboratory for Theoretical Computer Science

2. SEMANTICS OF BAYESIAN NETWORKS

- A Bayesian network for the random variables X_1, \dots, X_n is a representation of the joint probability distribution $\mathbf{P}(X_1, \dots, X_n)$.
- As before, a shorthand x_i is used for the atomic event $X_i = x_i$.
- Arrows encode conditional independence relations and therefore the probabilities of atomic events are determined by

$$P(x_1, \dots, x_n) = \prod_{i=1}^n P(x_i \mid \text{Parents}(x_i))$$

where $\text{Parents}(x_i)$ refers to the assignments of $Y \in \text{Parents}(X_i)$.

Example. Let us compute the probability of $j \wedge m \wedge a \wedge \neg b \wedge \neg e$:

$$\begin{aligned} &P(j \wedge m \wedge a \wedge \neg b \wedge \neg e) \\ = &P(j|a)P(m|a)P(a|\neg b \wedge \neg e)P(\neg e)P(\neg b) \\ = &0.9 \times 0.7 \times 0.001 \times 0.999 \times 0.998 = 0.00063. \end{aligned}$$

© 2006 TKK / Laboratory for Theoretical Computer Science

A Method for Constructing Bayesian Networks

- In a Bayesian network $G = \langle \{X_1, \dots, X_n\}, E \rangle$, a node $X_j \neq X_i$ is a *predecessor* of $X_i \iff$ there are nodes Y_1, \dots, Y_m such that $Y_1 = X_j$, $Y_m = X_i$, and $\forall j \in \{1, \dots, m-1\}: \langle Y_j, Y_{j+1} \rangle \in E$.
- Because G is a DAG, we may assume that the nodes X_1, \dots, X_n are ordered so that the predecessors of X_i are among X_1, \dots, X_{i-1} . Thus also $\text{Parents}(X_i) \subseteq \{X_1, \dots, X_{i-1}\}$.
- By the definition of conditional probability, we have that

$$\begin{aligned} P(x_1, \dots, x_n) &= \\ P(x_n \mid x_{n-1}, \dots, x_1) &P(x_{n-1}, \dots, x_1) = \\ P(x_n \mid x_{n-1}, \dots, x_1) &P(x_{n-1} \mid x_{n-2}, \dots, x_1) \cdots P(x_2 \mid x_1)P(x_1) = \\ \prod_{i=1}^n P(x_i \mid x_{i-1}, \dots, x_1). & \end{aligned}$$

© 2006 TKK / Laboratory for Theoretical Computer Science

- A Bayesian network is a correct representation if each variable X is conditionally independent of its predecessors Y given $\text{Parents}(X)$.
- Under the assumptions on conditional independence and node ordering, it can be established that

$$\mathbf{P}(X_i | X_{i-1}, \dots, X_1) = \mathbf{P}(X_i | \text{Parents}(X_i)). \quad (1)$$

- The choice of $\text{Parents}(X)$ for a random variable X affects how far conditional independence assumptions can be applied.
- $\text{Parents}(X)$ should contain all variables that directly influence X .

Example. Only *Alarm* directly influences *MaryCalls*. Given *Alarm*, *MaryCalls* is conditionally independent of the other variables:

$$\begin{aligned} \mathbf{P}(\text{MaryCalls} | \text{JohnCalls}, \text{Alarm}, \text{Earthquake}, \text{Burglary}) \\ = \mathbf{P}(\text{MaryCalls} | \text{Alarm}). \end{aligned}$$

Example. Let us reconstruct the Bayesian network for the alarm domain using a different node ordering:

MaryCalls, JohnCalls, Alarm, Burglary, Earthquake

1. As the first node, *MaryCalls* gets no parents.
2. When *JohnCalls* is added, *MaryCalls* becomes a parent of *JohnCalls*, as $\mathbf{P}(\text{JohnCalls} | \text{MaryCalls}) \neq \mathbf{P}(\text{JohnCalls})$.
3. Similarly, *Alarm* depends on both *MaryCalls* and *JohnCalls*.
4. Since $\mathbf{P}(\text{Burglary} | \text{Alarm}, \text{JohnCalls}, \text{MaryCalls}) = \mathbf{P}(\text{Burglary} | \text{Alarm})$, the only parent of *Burglary* is *Alarm*.
5. Nodes *Burglary* and *Alarm* become parents of *Earthquake*, as $\mathbf{P}(\text{Earthquake} | \text{Burglary}, \text{Alarm}, \text{JohnCalls}, \text{MaryCalls}) = \mathbf{P}(\text{Earthquake} | \text{Burglary}, \text{Alarm})$.

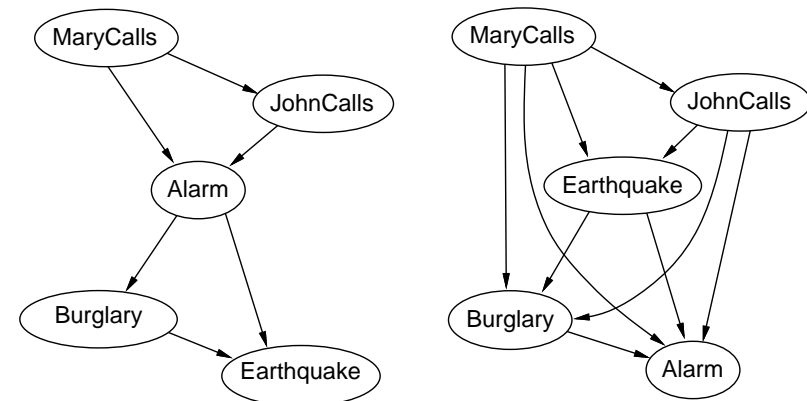
On Compactness and Node Ordering

- A Bayesian network can be a compact representation of the joint probability distribution (*locally structured* or *sparse* system).
- If each Boolean variable directly influences at most k other, then only $n2^k$ probabilities have to be specified (instead of 2^n).

Example. When $n = 30$ and $k = 5$, we would have to specify $n2^k = 960$ and $2^n = 1073741824$ probabilities, respectively.

- A clear **trade-off**: number of arrows (accuracy of probabilities) *versus* cost of specifying extra information (extending CPTs).
- Choosing a good node ordering is a non-trivial task.
- Heuristics: the *root causes* of the domain should be added first, then the variables influenced by them, and so forth.

- The resulting Bayesian network is given below on the left:

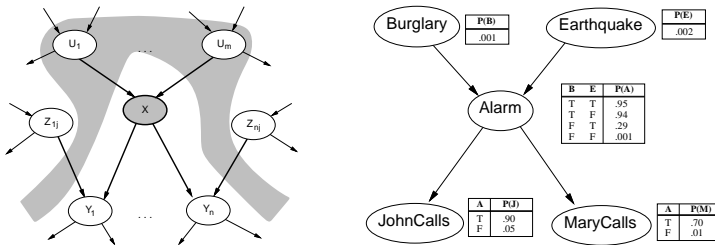


- The one on the right is obtained with another ordering and it is as complex (31 probabilities) as the full joint distribution!

Conditional Independence Relations

Two (equivalent) topological criteria can be utilized:

1. A node X is conditionally independent of its **non-descendants** (e.g., Z_{ij} s below), given its parents (i.e., U_i s below).

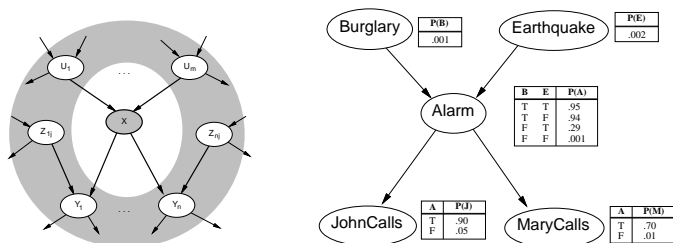


Example. In the burglary example, one may conclude that: *JohnCalls* is independent of *Burglary* and *Earthquake* given *Alarm*.

3. EFFICIENT REPRESENTATION OF CONDITIONAL DISTRIBUTIONS

- Specifying conditional probability tables means often a lot of work.
- To ease this process, some **canonical distributions** such as *deterministic* and *noisy logical* relationships have been proposed.
- When using a canonical distribution it is often enough to supply certain parameters rather than a complete CPT.
- There are also canonical continuous distributions such as *Gaussian distributions* and *probit/logit distributions*.

2. A node X is conditionally independent of all other nodes in the network, given its **Markov blanket** $mb(X)$, i.e., its parents, children, and children's parents.



Example. *Burglary* is independent of *JohnCalls* and *MaryCalls* given *Alarm* and *Earthquake*.

There is yet another criterion called **d-separation**, but unlike the first edition of the textbook it is not covered by the second.

Deterministic Nodes

- In the deterministic case, there is no uncertainty and the value of X is obtained as a (logical) function from those of $Parents(X)$.
- Deterministic nodes can also encode other fixed numerical functions depending on the variables involved.

Example. Define $NorthAmerican \leftrightarrow Canadian \vee US \vee Mexican$. This corresponds to specifying a CPT as follows:

<i>Canadian</i>	<i>US</i>	<i>Mexican</i>	<i>NorthAmerican</i>
F	F	F	0.0
T	F	F	1.0
⋮	⋮	⋮	⋮

Noisy Logical Relationships

- Noisy logical relationships add some uncertainty to the scenario.
- A **noisy OR** relationship comprises the following principles:
 1. Each cause has an independent chance of causing the effect.
 2. All possible causes are listed.
 3. Whatever inhibits some cause from causing an effect is independent of whatever inhibits other causes from causing the effect. Inhibitors are summarized as **noise parameters**.
- A noisy OR relationship in which a variable depends on k parents can be described using k parameters.
In contrast to this, 2^k entries are needed if a full CPT is specified.

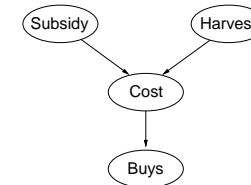
Bayesian networks with Continuous Variables

- Many real-world problems involve continuous quantities/variables.
- Continuous variables can be **discretized** but as a side-effect the resulting CPTs can become very large.
- Another possibility is to use standard probability density functions over the domains of continuous variables.
- A **hybrid Bayesian network** involves both discrete and continuous variables.

Example. Let us consider a medical domain including the variables *Fever* (a symptom), *Cold*, *Flu*, and *Malaria* (diseases). Using noise parameters $P(\neg fever \mid \neg cold, \neg flu, \neg malaria) = 0.6$, $P(\neg fever \mid \neg cold, flu, \neg malaria) = 0.2$, and $P(\neg fever \mid \neg cold, \neg flu, malaria) = 0.1$, we get the following CPT:

<i>Cold</i>	<i>Flu</i>	<i>Malaria</i>	$P(\text{Fever})$	$P(\neg \text{Fever})$
<i>F</i>	<i>F</i>	<i>F</i>	0.0	1.0
<i>F</i>	<i>F</i>	<i>T</i>	0.9	0.1
<i>F</i>	<i>T</i>	<i>F</i>	0.8	0.2
<i>F</i>	<i>T</i>	<i>T</i>	0.98	$0.02 = 0.2 \times 0.1$
<i>T</i>	<i>F</i>	<i>F</i>	0.4	0.6
<i>T</i>	<i>F</i>	<i>T</i>	0.94	$0.06 = 0.6 \times 0.1$
<i>T</i>	<i>T</i>	<i>F</i>	0.88	$0.12 = 0.6 \times 0.2$
<i>T</i>	<i>T</i>	<i>T</i>	0.988	$0.012 = 0.6 \times 0.2 \times 0.1$

Example. Consider a system one Boolean random variable *Subsidy* and three continuous random variables *Harvest*, *Cost*, and *Buys*.



For *Cost*, we need to specify $\mathbf{P}(\text{Cost} \mid \text{Harvest}, \text{Subsidy})$.

- The discrete parent is handled by explicitly enumerating both $\mathbf{P}(\text{Cost} \mid \text{Harvest}, \text{subsidy})$ and $\mathbf{P}(\text{Cost} \mid \text{Harvest}, \neg \text{subsidy})$.
- The parameters of the cost distribution (e.g. **linear Gaussian** distribution) are given as a function of the variable *Harvest*.

4. EXACT INFERENCE IN BAYESIAN NETWORKS

- An agent gets values for evidence variables from its percepts and asks about the possible values of other variables so that it can decide what action to take (recall the decision theoretic design).
- The basic task of a probabilistic reasoning system is to compute $\mathbf{P}(X | E_1 = e_1, \dots, E_m = e_m)$ given a **query variable** X and exact values e_1, \dots, e_m of some **evidence variables** E_1, \dots, E_m .
- The remaining variables Y_1, \dots, Y_n act as **hidden variables**.

Examples. Recalling the alarm example, the problem is to calculate distributions such as $\mathbf{P}(\text{Burglary} | \text{JohnCalls}, \text{MaryCalls})$ and $\mathbf{P}(\text{Alarm} | \text{JohnCalls}, \text{Earthquake})$?

Example. Consider the query $\mathbf{P}(B | j, m)$ in the burglary example.

For this query, E and A are hidden variables and enumeration amounts to computing the following distribution (in a depth first fashion):

$$\begin{aligned}
 \mathbf{P}(B | j, m) &= \alpha \mathbf{P}(B, j, m) \\
 &= \alpha \sum_e \sum_a \mathbf{P}(B, e, a, j, m) \\
 &= \alpha \sum_e \sum_a \mathbf{P}(B) P(e) \mathbf{P}(a | B, e) P(j | a) P(m | a) \\
 &= \alpha \mathbf{P}(B) \sum_e P(e) \sum_a \mathbf{P}(a | B, e) P(j | a) P(m | a) \\
 &= \alpha \langle 0.00059224, 0.0014919 \rangle \\
 &\approx \langle 0.284, 0.716 \rangle
 \end{aligned}$$

The details of computing $P(b | j, m)$ are analyzed next.

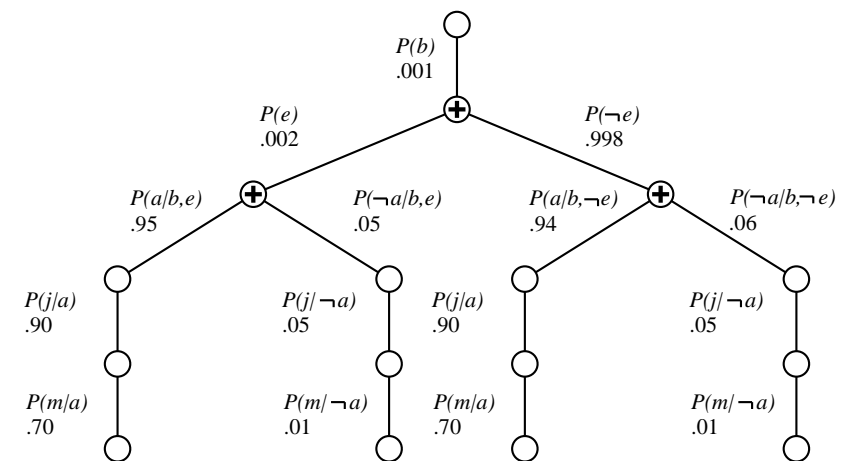
Inference by Enumeration

- We introduce shorthands \mathbf{E} and \mathbf{Y} for E_1, \dots, E_m and Y_1, \dots, Y_n , respectively, and similarly \mathbf{e} and \mathbf{y} for their values.
- A query $\mathbf{P}(X | \mathbf{e})$ can be answered by exhaustive enumeration:

$$\mathbf{P}(X | \mathbf{e}) = \alpha \mathbf{P}(X, \mathbf{e}) = \alpha \sum_{\mathbf{y}} \mathbf{P}(X, \mathbf{e}, \mathbf{y})$$

where α is a normalizing constant.

- If a Bayesian network is used, this leads to the computation of sums of products of conditional probabilities from the network.
- The time complexity for a network of n variables is of order 2^n .



➡ Certain subexpressions are computed repeatedly.

Variable Elimination Algorithm

- The enumeration algorithm can be improved substantially by doing calculations in a bottom-up fashion using **factors** which are matrices of probabilities.
- The **pointwise product** of two factors $f_1(\mathbf{X}, \mathbf{Y})$ and $f_2(\mathbf{Y}, \mathbf{Z})$ is defined by $(f_1 \times f_2)(\mathbf{X}, \mathbf{Y}, \mathbf{Z}) = f_1(\mathbf{X}, \mathbf{Y})f_2(\mathbf{Y}, \mathbf{Z})$.
- A variable X can be **summed out** from a product of factors $f_i(X, \mathbf{Y})$ by computing $\sum_x (f_1(x, \mathbf{Y}) \times \dots \times f_n(x, \mathbf{Y}))$.
- Multiplication takes place only when summing out variables.
- Every variable that is not an ancestor of a query variable or evidence variable is irrelevant to the query and thus removable.

The Complexity of Exact Inference

- A **polytree** is a **singly connected** graph: there is at most one undirected path between any two nodes.
- If a belief network forms a polytree, the probability distribution $\mathbf{P}(X | e)$ can be computed very efficiently (in **linear time**).
- For **multiply connected** networks, in which at least two variables are connected by several paths, variable elimination can have exponential time and space complexity in the worst case.
- In general, exact inference in Bayesian networks is NP-hard (even #P-hard) as it includes propositional inference as a special case.

Example. The computation of the previous distribution

$$\mathbf{P}(B | j, m) = \alpha \mathbf{P}(B) \sum_e P(e) \sum_a \mathbf{P}(a | B, e) P(j | a) P(m | a)$$

takes place bottom-up using factors as follows:

1. $f_M(A) = \langle P(m | a), P(m | \neg a) \rangle$;
2. $f_J(A)$ is defined analogously;
3. $f_A(A, B, E)$ is three-dimensional;
4. the variable A is summed out from the product of these three:

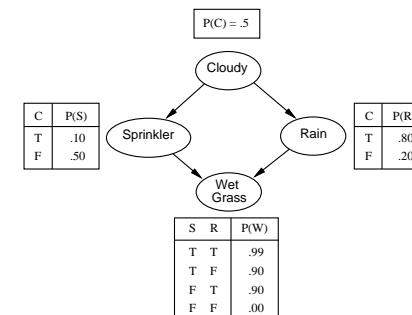
$$f_{J,M}(B, E) = \sum_a (f_A(a, B, E) \times f_J(a) \times f_M(a));$$

5. E is summed out similarly and $\mathbf{P}(B | j, m) = \alpha f_B(B) \times f_{J,M}(B)$.

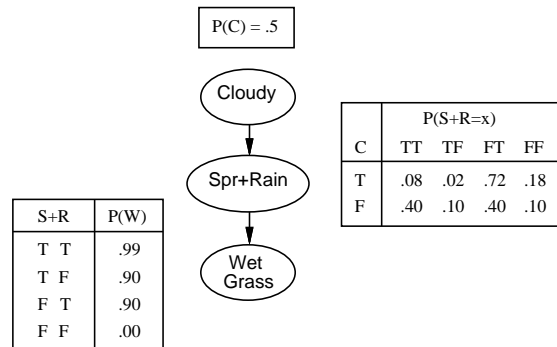
Clustering Algorithms

- Multiply connected Bayesian networks can be transformed into polytrees by combining some nodes into cluster nodes.

Example. Consider clustering the nodes *Sprinkler* and *Rain* in the following multiply connected network:



- The following polytree network is obtained:



- Linear time algorithms can be used for query answering, but the size of the network grows exponentially in the worst case.
- Typically, there are several ways to compose cluster nodes and it is non-trivial to choose the best way to perform clustering.

Direct Sampling Methods

- In direct sampling, the world described by a Bayesian network (without evidence) is simulated stochastically.
- Atomic events are randomly generated in topological order by selecting definite values for random variables.
- The value for a random variable X is chosen according to the conditional probability table associated with X .
- **Prior sampling** produces the event x_1, \dots, x_n with probability

$$S_{PS}(x_1, \dots, x_n) = \prod_{i=1}^n P(x_i | \text{Parents}(X_i)) = P(x_1, \dots, x_n).$$

5. APPROXIMATE INFERENCE IN BAYESIAN NETWORKS

- Randomized sampling algorithms provide approximate answers whose accuracy depends on the number of samples generated.
- Here sampling is applied to the computation of posterior probabilities given a prior distribution (a Bayesian network).
- There are several approximation methods including
 - Direct sampling
 - Rejection sampling
 - Likelihood weighting
 - Markov chain Monte Carlo algorithm

- The posterior distribution $\mathbf{P}(X | \mathbf{e}) = \frac{\mathbf{P}(X, \mathbf{e})}{\mathbf{P}(\mathbf{e})}$ is estimated by counting the frequencies with which events occur.

- The number of samples N affects accuracy:

$$\lim_{N \rightarrow \infty} \frac{N_{PS}(x_1, \dots, x_n)}{N} = S_{PS}(x_1, \dots, x_n) = P(x_1, \dots, x_n).$$

- Direct sampling is not very useful if the event \mathbf{e} occurs very rarely.

Example. Let us produce one sample the lawn watering domain:

$$\mathbf{P}(\text{Cloudy}) = \langle 0.5, 0.5 \rangle \quad \Rightarrow \quad \text{return true}$$

$$\mathbf{P}(\text{Sprinkler} | \text{cloudy}) = \langle 0.1, 0.9 \rangle \quad \Rightarrow \quad \text{return false}$$

$$\mathbf{P}(\text{Rain} | \text{cloudy}) = \langle 0.8, 0.2 \rangle \quad \Rightarrow \quad \text{return true}$$

$$\mathbf{P}(\text{WetGrass} | \neg \text{sprinkler}, \text{rain}) = \langle 0.9, 0.1 \rangle \quad \Rightarrow \quad \text{return true}$$

Example. E.g., $\mathbf{P}(\text{WetGrass} | \text{sprinkler} \wedge \text{rain})$ converges slowly.

Rejection Sampling in Bayesian Networks

- In its simplest form, rejection sampling can be used to compute conditional probabilities such as $P(X | \mathbf{e})$.
- Samples are generated from the prior distribution, but samples which do not match the evidence are rejected.
- The estimated distribution $\hat{\mathbf{P}}(X | \mathbf{e}) = \alpha N_{PS}(X, \mathbf{e}) = \frac{N_{PS}(X, \mathbf{e})}{N_{PS}(\mathbf{e})}$.
- With sufficiently many samples $\hat{\mathbf{P}}(X | \mathbf{e}) \approx \frac{P(X, \mathbf{e})}{P(\mathbf{e})} = \mathbf{P}(X | \mathbf{e})$.
- Rejection sampling tends to reject too many samples.

Example. Suppose that out of 100 samples, 73 are rejected as *Sprinkler = false*. Out of the remaining 27 samples, 8 satisfy *Rain = true*. Thus $\mathbf{P}(Rain | sprinkler) \approx \alpha \langle 8, 19 \rangle = \langle 0.296, 0.704 \rangle$.

Example. Let us estimate $\mathbf{P}(Rain | sprinkler, wetgrass)$ by likelihood weighting. Initially, the weight w is set to 1.0.

The values of variables are chosen randomly as follows:

1. $\mathbf{P}(Cloudy) = \langle 0.5, 0.5 \rangle \implies cloudy$ is randomly chosen.
2. *Sprinkler* is an evidence variable that has been set to *true*:
 w is revised to $w \times P(sprinkler | cloudy) = 0.1$.
3. $\mathbf{P}(Rain | cloudy) = \langle 0.8, 0.2 \rangle \implies rain$ is randomly chosen.
4. *WetGrass* is an evidence variable with value *true*:
 w is revised to $w \times P(wetgrass | sprinkler, rain) = 0.099$.

☞ We have completed a run saying that *Rain = true* given *sprinkler* and *wetgrass* with a likelihood weight 0.099.

Likelihood Weighting

- Likelihood weighting is similar to rejection sampling, but the values of evidence variables \mathbf{E} are kept fixed while sampling others.
- The CPTs of the Bayesian network are consulted to see how likely the event \mathbf{e} is.
- In this way, the conditional probability $P(\mathbf{e} | x, \mathbf{y})$ is interpreted as a likelihood weight for that particular run.
- An estimate of $P(X = x | \mathbf{e})$ is obtained as a weighted proportion of runs with $X = x$ among the runs accumulated so far.
- Likelihood weighting converges faster than rejection sampling.
- Getting accurate probabilities for unlikely events is still a problem.

Why Likelihood Weighting Works

- The algorithm samples each non-evidence variable in $\mathbf{Z} = \{X\} \cup \mathbf{Y}$ given the values of its parents:

$$S_{WS}(\mathbf{z}, \mathbf{e}) = \prod_{i=1}^l P(z_i | \text{Parents}(Z_i)).$$

- The weight for a given sample is $w(\mathbf{z}, \mathbf{e}) = \prod_{i=1}^m P(e_i | \text{Parents}(E_i))$.
- The *weighted* probability $S_{WS}(\mathbf{z}, \mathbf{e})w(\mathbf{z}, \mathbf{e}) = P(\mathbf{y}, \mathbf{e})$.
- Likelihood weighting estimates are shown consistent as follows:

$$\begin{aligned} \hat{P}(x | \mathbf{e}) &= \alpha \sum_{\mathbf{y}} N_{WS}(x, \mathbf{y}, \mathbf{e}) w(x, \mathbf{y}, \mathbf{e}) \\ &\approx \alpha' \sum_{\mathbf{y}} S_{WS}(x, \mathbf{y}, \mathbf{e}) w(x, \mathbf{y}, \mathbf{e}) \quad (\text{for large } N) \\ &= \alpha' \sum_{\mathbf{y}} P(x, \mathbf{y}, \mathbf{e}) \\ &= \alpha' P(x, \mathbf{e}) = P(x | \mathbf{e}). \end{aligned}$$

Inference by Markov Chain Simulation

- A **Markov chain Monte Carlo** (MCMC) algorithm generates the next state by sampling a value for a nonevidence variable X_i conditioned by the current values of the variables in $mb(X_i)$.
- The simulation starts from a random state \mathbf{x} for $\mathbf{X} = \{X\} \cup \mathbf{Z}$.
- Each round of the simulation consists of the following steps:
 1. Increase the count $\mathbf{N}[x]$ by one for the current value x of X .
 2. Sample the value of each X_i in \mathbf{X} using $\mathbf{P}(X_i | mb(X_i))$.
- The estimate for the distribution $\mathbf{P}(X | \mathbf{e})$ is obtained by normalizing the counts in $\mathbf{N}[X]$.

Detailed Balance

- One interpretation of the equation $\pi(\mathbf{x}') = \sum_{\mathbf{x}} \pi(\mathbf{x})q(\mathbf{x} \rightarrow \mathbf{x}')$ is that the expected *outflow* from each state (population) is equal to the expected *inflow* from all the states.
- Assuming the equality of flows in both directions leads to the property of **detailed balance**: for all \mathbf{x} and \mathbf{x}' :

$$\pi(\mathbf{x})q(\mathbf{x} \rightarrow \mathbf{x}') = \pi(\mathbf{x}')q(\mathbf{x}' \rightarrow \mathbf{x}).$$

- Stationarity is implied by detailed balance:

$$\begin{aligned} \sum_{\mathbf{x}} \pi(\mathbf{x})q(\mathbf{x} \rightarrow \mathbf{x}') &= \sum_{\mathbf{x}} \pi(\mathbf{x}')q(\mathbf{x}' \rightarrow \mathbf{x}) \\ &= \pi(\mathbf{x}') \sum_{\mathbf{x}} q(\mathbf{x}' \rightarrow \mathbf{x}) \\ &= \pi(\mathbf{x}'). \end{aligned}$$

Why MCMC Works

- The sampling process settles into a “dynamic equilibrium” in which the long-run fraction of time spent in each state is exactly proportional to its posterior probability.
- A **Markov chain** is defined by **transition probabilities** $q(\mathbf{x} \rightarrow \mathbf{x}')$ from a state \mathbf{x} to a state \mathbf{x}' .
- Let $\pi_t(\mathbf{x})$ denote the probability of a state \mathbf{x} after t steps.
- For the next step, we have $\pi_{t+1}(\mathbf{x}') = \sum_{\mathbf{x}} \pi_t(\mathbf{x})q(\mathbf{x} \rightarrow \mathbf{x}')$.

Definition. The chain has reached its **stationary distribution** π if $\pi_{t+1} = \pi_t$, i.e., π is defined by $\pi(\mathbf{x}') = \sum_{\mathbf{x}} \pi(\mathbf{x})q(\mathbf{x} \rightarrow \mathbf{x}')$ for all \mathbf{x}' .

Gibbs sampler

- Let X_i be the variable to be sampled and \bar{X}_i all the hidden variables other than X_i .
- The **Gibbs sampler** is based on transition probabilities

$$q(\mathbf{x} \rightarrow \mathbf{x}') = q((x_i, \bar{x}_i) \rightarrow (x'_i, \bar{x}_i)) = P(x'_i | \bar{x}_i, \mathbf{e}) = P(x'_i | mb(X_i)).$$
- Gibbs sampler is in detailed balance with the true posterior:

$$\begin{aligned} \pi(\mathbf{x})q(\mathbf{x} \rightarrow \mathbf{x}') &= P(\mathbf{x} | \mathbf{e})P(x'_i | \bar{x}_i, \mathbf{e}) \\ &= P(x_i, \bar{x}_i, \mathbf{e})P(x'_i | \bar{x}_i, \mathbf{e}) \\ &= P(x_i | \bar{x}_i, \mathbf{e})P(\bar{x}_i | \mathbf{e})P(x'_i | \bar{x}_i, \mathbf{e}) \\ &= P(x'_i, \bar{x}_i, \mathbf{e})P(x_i | \bar{x}_i, \mathbf{e}) \\ &= \pi(\mathbf{x}')q(\mathbf{x}' \rightarrow \mathbf{x}). \end{aligned}$$

6. OTHER APPROACHES TO UNCERTAIN REASONING

- Early expert systems were based on strict logical reasoning.
- Probabilistic techniques were dominating in the second generation, but these techniques suffered from the exponential blow-up of the joint probability distribution w.r.t. the number of variables.
- Consequently, many alternatives to probabilities were pursued:
 1. Default reasoning
 2. Rules with certainty factors
 3. Dempster-Shafer theory
 4. Fuzzy logic

- Logic programs with *negation as failure to prove* from an important subclass of non-monotonic theories.

Example. Let us describe the applicability of actions using rules:

$$\{ \begin{array}{l} doable(A) \leftarrow preconds(A) \wedge \text{not } exceptional(A), \\ exceptional(A) \leftarrow \text{not } deterministic(A), \\ exceptional(A) \leftarrow delayed(A) \end{array} \}$$

- The semantics of “not ϕ ” is different from classical negation $\neg\phi$.
- The conclusion $doable(run)$ can be drawn by the rules above given the premises $preconds(run)$ and $deterministic(run)$.
- Such a conclusion is no longer possible if $delayed(run)$ is introduced as an additional premise.
- Dropping the premise $deterministic(A)$ has the same effect.

Default Reasoning

- Reasoning by default means inferring something in the absence of any information to the contrary.
- Provides a compact way to encode exceptions to general principles.
- A qualitative approach to handle uncertainty.
- Default reasoning *violates* the **monotonicity** property of classical logic: if $\Sigma_1 \models \phi$ and $\Sigma_1 \subseteq \Sigma_2$, then $\Sigma_2 \models \phi$.
- Several formalizations of *non-monotonic reasoning* have been proposed: **default logic** [Reiter, 1980], **circumscription** [McCarthy, 1980], **autoepistemic logic** [Moore, 1983], ...
- Implementation techniques have substantially improved during 90s.

Logical Rules and Certainty Factors

- Reasoning systems based on classical logic have important properties that are lacked by their probabilistic counterparts:
 1. **Locality:** a rule can be used for making inferences without worrying about the other rules in the system.
 2. **Detachment:** if a sentence ϕ is proven to be valid, it can be detached from its justification (proof), as it universally true.
 3. **Truth-functionality:** the truth values of complex sentences can be computed from the truth values of their components.
- Unfortunately, problems arise with truth-functionality and chained inferences, if logical rules are equipped with certainty factors.

Example. For instance, $Sprinkler \mapsto WetGrass$ and $WetGrass \mapsto Rain$ tend to imply $Sprinkler \mapsto Rain$.

Dempster-Shafer theory

- Dempster-Shafer theory has been designed to deal with the distinction between **uncertainty** and **ignorance**.
- The *belief function* $Bel(X)$ gives the probability that the evidence obtained so far supports X .

Example. Consider flipping a coin under the following circumstances:

1. If the coin is doubted to be unfair (nothing can be assumed about its behavior), then $Bel(heads) = 0$ and $Bel(\neg heads) = 0$.
2. If the coin is fair with a certainty of 0.9, then we have $Bel(heads) = 0.5 \times 0.9 = 0.45$ and $Bel(\neg heads) = 0.45$

☞ We obtain *probability intervals* $[0, 1]$ and $[0.45, 0.55]$ for *Heads*.

SUMMARY

- **Conditional independence** information can be used for structuring and simplifying knowledge about an uncertain domain.
- **Bayesian networks** provide a natural way to represent conditional independence information.
- A Bayesian network is a complete (and often also very compact) representation of the joint probability distribution.
- Efficient algorithms exist for Bayesian networks that are topologically *polytrees*, but reasoning with Bayesian networks is NP-hard in general.
- Probabilities can be estimated by **sampling methods**.

Fuzzy Logic

- **Fuzzy set theory** is about specifying how well an object satisfies a vague description rather than uncertainty.

Example. For instance, a statement like "Mika Myllylä is tall" can be assigned a truth value between 0 and 1 (even if it is known how tall he is).

- The *fuzzy truth* of complex sentences is defined truth-functionally:

$$T(\phi \wedge \psi) = \min(T(\phi), T(\psi)),$$

$$T(\phi \vee \psi) = \max(T(\phi), T(\psi)), \text{ and}$$

$$T(\neg A) = 1 - T(A).$$

- Despite of semantic difficulties, fuzzy logic has been very successful in commercial applications involving *automated control*.

QUESTIONS

- Build a Bayesian network for the soccer domain.
 1. Choose appropriate variables for the description of the domain.
 2. Choose an ordering for the variables.
 3. Construct the actual belief network by
 - (i) analyzing dependencies among variables and
 - (ii) defining CPTs for each variable.
- Make both causal and diagnostic inferences using the network.