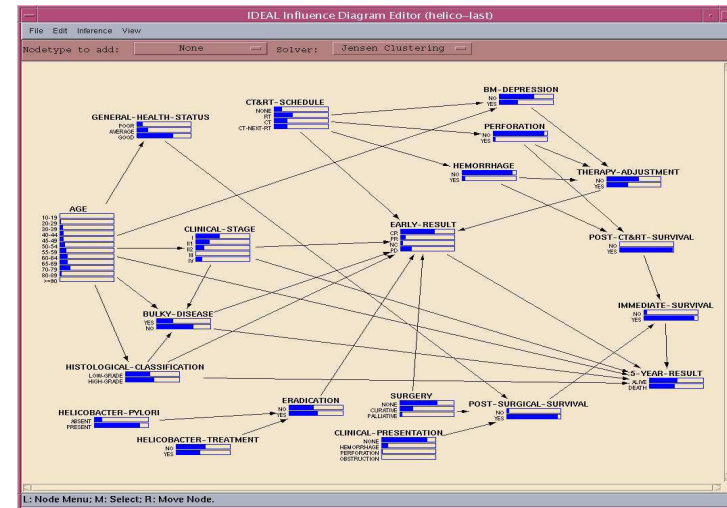


Example of a Bayesian network



Bayesian networks

Principles and Definitions

Lecture2: Bayesian networks – p. 1

Why Bayesian networks?

Probabilistic graphical models, such as Bayesian networks, are now the most popular uncertainty formalisms because:

- Handle noise, missing information and probabilistic relations
- Learn from data and can incorporate domain knowledge
- Offer flexible reasoning
- Have compact graphical representation (interface)
- Founded principles: probability theory
- Engineering principles: knowledge acquisition, machine learning and statistics

Lecture2: Bayesian networks – p. 3

Lecture2: Bayesian networks – p. 2

Popular applications

- Hardware trouble shooting: Microsoft, Boeing, HP
- Biological modelling: gene expressions
- Medical diagnosis and therapy selection: BNs are now the most popular paradigm for medical intelligent systems
- Art: orchestral music accompaniment
- and more ...

Lecture2: Bayesian networks – p. 4

General notation

- **Stochastic (= statistical = random) variable:** upper-case letter, e.g. X , or upper-case string, e.g. **FEVER**
- **Values:** variables can take on values, e.g. $X = x$, **FEVER = yes**
- **Binary variables:** take one of *two* values, e.g. $X = true$ and $X = false$
- **Discrete variables:** take only one of a finite set of possible values, e.g. $TEMP = \{low, medium, high\}$
- **Continuous variables:** take any value in some interval or intervals of real numbers \mathcal{R} , e.g. $TEMP \in [-50, 50]$

Lecture2: Bayesian networks – p. 5

Abbreviated notation

- **Binary variables:** $X = true$ as x , and $X = false$ as $\neg x$
- **Non-binary variables:** $X = x$ as x or $CITY = tokyo$ as $tokyo$
- **Sets of variables:** analogous to variables
 - Example:

$$\begin{array}{l}
 X_1 = x_1 \\
 X_2 = x_2 \\
 \cdot \\
 \cdot \\
 X_n = x_n
 \end{array}
 \implies
 \begin{array}{l}
 X = (X_1, X_2, \dots, X_n) \\
 x = (x_1, x_2, \dots, x_n) \\
 X = x
 \end{array}$$

Lecture2: Bayesian networks – p. 6

Abbreviated notation (cont.)

- **Conjunctions:** $(X = x) \wedge (Y = y)$ as $(X = x, Y = y)$
- **Templates:** (X, Y) means $(X = x, Y = y)$, for *any* value x, y , i.e. the choice of the values x and y does not really matter
- **Examples:**
 - $P(X = x, Y = y) \Leftrightarrow P(X = x \wedge Y = y)$
 - $P(X, Y) \Leftrightarrow P(X = x, Y = y)$, for *any* value x, y
 - $P(X | Y) \Leftrightarrow P(X = x | Y = y)$, for *any* value x, y
- $\sum_X P(X) = P(x) + P(\neg x)$, where X is binary

Lecture2: Bayesian networks – p. 7

Probability theory

- **Probability distribution P :** attaches a number in (closed) interval $[0, 1]$ to *Boolean expressions*
- **Boolean algebra \mathbb{B}** (for two variables RAIN and HAPPY):
 - \top (*true*),
rain, $\neg rain$,
happy, $\neg happy$,
rain \wedge *happy*, ..., *rain* \wedge *happy* \wedge $\neg happy$, ...,
 $\neg rain$ \wedge *happy*, ..., *rain* \vee *happy*,
 - \perp (*false*)
- such that:
 - $\perp \leq rain$, $rain \leq (rain \vee happy)$, ... (in general $\perp \leq x$ for each Boolean expression $x \in \mathbb{B}$);
 - $x \leq \top$ for each Boolean expression $x \in \mathbb{B}$

Lecture2: Bayesian networks – p. 8

Probability distribution

- A **probability distribution** P is defined as a function $P : \mathbb{B} \rightarrow [0, 1]$, such that:
 - $P(\perp) = 0$
 - $P(\top) = 1$
 - $P(x \vee y) = P(x) + P(y)$, if $x \wedge y = \perp$ with $x, y \in \mathbb{B}$
- Examples:
 - $P(\text{rain} \vee \text{happy}) = P(\text{rain}) + P(\text{happy})$, as $\text{rain} \wedge \text{happy} = \perp$ (why? Because I define it that way)
 - $P(\text{rain} \wedge \text{happy}) = P(\perp) = 0$
 - $P(\neg \text{rain} \vee \text{rain}) = P(\neg \text{rain}) + P(\text{rain}) = P(\top) = 1 \Rightarrow P(\neg \text{rain}) = 1 - P(\text{rain})$
 - $0 \leq P(\text{rain}) \leq 1$

Lecture2: Bayesian networks – p. 9

Joint probability distribution

Let X and Y be random variables with $\{X\} = \{x_1, x_2, \dots, x_n\}$ and $\{Y\} = \{y_1, y_2, \dots, y_m\}$.

The product set

$$\{X\} \times \{Y\} = \{x_1, x_2, \dots, x_n\} \times \{y_1, y_2, \dots, y_m\}$$

is made into a probability space by defining

$$P(X = x_i \wedge Y = y_j) = P(x_i, y_j)$$

where P is a **joint probability function**

Lecture2: Bayesian networks – p. 11

Probability distribution (cont.)

- **Boolean algebras** \Leftrightarrow **sets**:
 - $\top \Leftrightarrow \Omega$
 - $\perp \Leftrightarrow \emptyset$
 - $x \Leftrightarrow X$
 - $\neg x \Leftrightarrow \bar{X}$
 - $(x \vee y) \Leftrightarrow (X \cup Y)$
 - $(x \wedge y) \Leftrightarrow (X \cap Y)$
 - $x \leq (x \vee y) \Leftrightarrow X \subseteq (X \cup Y)$

with \Leftrightarrow 1-1 correspondence, e.g.

$$P(\overline{\text{Rain}}) = 1 - P(\text{Rain})$$

Lecture2: Bayesian networks – p. 10

Marginalisation

Suppose the joint probability distribution of two variables X and Y is given; then

$$\begin{aligned} P(x) = P(X = x) &= P(x \wedge \top) \\ &= P(x \wedge (y \vee \bar{y})) \\ &= P((x \wedge y) \vee (x \wedge \bar{y})) \\ &= P(x \wedge y) + P(x \wedge \bar{y}) \end{aligned}$$

since $P(a \vee b) = P(a) + P(b)$, if $a \wedge b = \perp$

$$\Rightarrow P(x) = \sum_Y P(x, Y)$$

also known as **marginal probability function** of X .

Lecture2: Bayesian networks – p. 12

Example

- Assume that X_1, X_2, X_3 and X_4 are binary variables. Then $P(X_1, X_2, X_3, X_4)$:

$P(x_1, x_2, x_3, x_4) = 0.1$	$P(x_1, \neg x_2, \neg x_3, x_4) = 0.015$
$P(x_1, \neg x_2, x_3, x_4) = 0.04$	$P(x_1, \neg x_2, x_3, \neg x_4) = 0.1$
$P(x_1, x_2, \neg x_3, x_4) = 0.03$	$P(x_1, x_2, \neg x_3, \neg x_4) = 0.004$
$P(x_1, x_2, x_3, \neg x_4) = 0.1$	$P(\neg x_1, \neg x_2, \neg x_3, x_4) = 0.005$
$P(\neg x_1, x_2, x_3, x_4) = 0.0$	$P(\neg x_1, \neg x_2, x_3, \neg x_4) = 0.01$
$P(\neg x_1, \neg x_2, x_3, x_4) = 0.2$	$P(\neg x_1, x_2, \neg x_3, \neg x_4) = 0.01$
$P(\neg x_1, x_2, \neg x_3, x_4) = 0.08$	$P(x_1, \neg x_2, \neg x_3, \neg x_4) = 0.006$
$P(\neg x_1, x_2, x_3, \neg x_4) = 0.1$	$P(\neg x_1, \neg x_2, \neg x_3, \neg x_4) = 0.2$

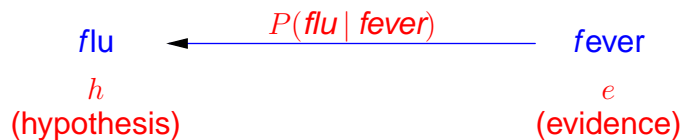
- $\sum_{X_1, X_2, X_3, X_4} P(X_1, X_2, X_3, X_4) = 1$
- Marginalisation:

$$P(x_4) = \sum_{X_1, X_2, X_3} P(X_1, X_2, X_3, x_4) = 0.47$$

Lecture2: Bayesian networks – p. 13

Reversal of chances

- $P(\text{flu} | \text{fever})$ is usually **unknown**:

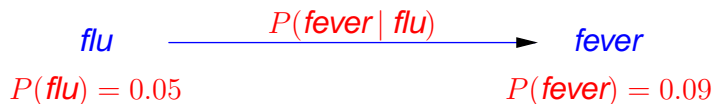


- Known is:

$$P(\text{fever} | \text{flu}) = 0.9$$

$$P(\text{flu}) = 0.05$$

$$P(\text{fever}) = 0.09$$



Lecture2: Bayesian networks – p. 15

Conditional probability

(Example: *flu* and *fever*)

- $P(\text{flu} \wedge \text{fever})$: chance of *flu* and *fever* at the same time
- $P(\text{flu} | \text{fever})$: chance of *flu* knowing that the person already has *fever* (conditional probability)
- Definition:

$$P(\text{flu} | \text{fever}) = \frac{P(\text{flu} \wedge \text{fever})}{P(\text{fever})}$$

↗
adjust $P(\text{flu} \wedge \text{fever})$, so that uncertainty in 'fever' is removed

Lecture2: Bayesian networks – p. 14

Bayes' rule

"I now send you an essay which I have found among the papers of our deceased friend Mr Bayes, and which, in my opinion, has great merit... In an introduction which he has writ to this Essay, he says, that his design at first in thinking on the subject of it was, to find out **a method by which we might judge concerning the probability that an event has to happen, in given circumstances, upon supposition that we know nothing concerning it but that, under the same circumstances, it has happened a certain number of times, and failed a certain other number of times.**"

Richard Price

Introducing "Essay towards solving a problem in the doctrine of chances" by Thomas Bayes to the Royal Society of London in 1764

Lecture2: Bayesian networks – p. 16

Bayes' rule - Example

- Bayes' rule – reversal of chances:

$$P(e | h) \quad P(\text{fever} | \text{flu}) = 0.9$$

$$P(h) \quad P(\text{flu}) = 0.05$$

$$P(e) \quad P(\text{fever}) = 0.09$$

$$\begin{aligned} P(\text{flu} | \text{fever}) &= \frac{P(\text{fever} | \text{flu})P(\text{flu})}{P(\text{fever})} \\ &= 0.9 \cdot 0.05 / 0.09 = 0.5 \end{aligned}$$

- Definition of Bayes' rule (the 'chance reverter'):

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

Lecture2: Bayesian networks – p. 17

Chain rule (definition)

$$\begin{aligned} P(X_1, X_2, \dots, X_n) &= P(X_1 | X_2, \dots, X_n) \cdot \\ &\quad P(X_2 | X_3, \dots, X_n) \cdot \\ &\quad P(X_3 | X_4, \dots, X_n) \cdot \\ &\quad \vdots \\ &\quad P(X_{n-1} | X_n) \cdot \\ &\quad P(X_n) \\ &= \prod_{i=1}^{n-1} P(X_i | X_{i+1}, \dots, X_n) P(X_n) \end{aligned}$$

Lecture2: Bayesian networks – p. 19

Chain rule (derivation)

Definition of conditional probability:

$$P(X_1 | X_2, \dots, X_n) = \frac{P(X_1, X_2, \dots, X_n)}{P(X_2, \dots, X_n)}$$

$$\Rightarrow P(X_1, X_2, \dots, X_n) = P(X_1 | X_2, \dots, X_n) P(X_2, \dots, X_n)$$

Furthermore,

$$P(X_2, \dots, X_n) = P(X_2 | X_3, \dots, X_n) P(X_3, \dots, X_n)$$

\vdots

$$P(X_{n-1}, X_n) = P(X_{n-1} | X_n) P(X_n)$$

$$P(X_n) = P(X_n)$$

Lecture2: Bayesian networks – p. 18

Definition Bayesian network (BN)

A Bayesian network \mathcal{B} is a pair $\mathcal{B} = (G, P)$, where:

- $G = (V(G), A(G))$ is an **acyclic directed graph**, with
 - $V(G) = \{v_1, v_2, \dots, v_n\}$, a set of **vertices** (nodes)
 - $A(G) \subseteq V(G) \times V(G)$ a set of **arcs**
- $P : \wp(V(G)) \rightarrow [0, 1]$ is a **joint probability distribution**, such that

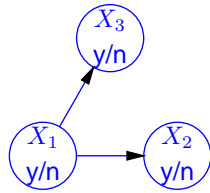
$$P(V(G)) = \prod_{i=1}^n P(v_i | \pi_G(v_i))$$

where $\pi_G(v_i)$ denotes the set of immediate ancestors (parents) of vertex v_i in G

- Notational convenience: $v_i \longrightarrow X_i$

Lecture2: Bayesian networks – p. 20

Example of a Bayesian network



Bayesian network $\mathcal{B} = (G, P)$, where $G = (V(G), A(G))$, with

- Set of vertices: $V(G) = \{X_1, X_2, X_3\}$
- Set of arcs: $A(G) = \{(X_1, X_2), (X_1, X_3)\}$
- Joint probability distribution:

$$P(X_1, X_2, X_3) = P(X_1) \cdot P(X_2 | X_1) \cdot P(X_3 | X_1)$$

Lecture2: Bayesian networks – p. 21

Example (cont.)

$$P(X_1, X_2, X_3) = P(X_1) \cdot P(X_2 | X_1) \cdot P(X_3 | X_1)$$

with for example:

$$P(x_1) = 0.7$$

$$P(\neg x_1) = 0.3 = 1 - P(x_1)$$

$$P(x_2 | x_1) = 0.6$$

$$P(\neg x_2 | x_1) = 0.4$$

$$P(x_2 | \neg x_1) = 0.1$$

$$P(\neg x_2 | \neg x_1) = 0.9$$

$$P(x_3 | x_1) = 0.1$$

$$P(\neg x_3 | x_1) = 0.9$$

$$P(x_3 | \neg x_1) = 0.8$$

$$P(\neg x_3 | \neg x_1) = 0.2$$

Lecture2: Bayesian networks – p. 22

Conditional independence relation

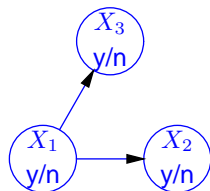
Let X, Y, Z be sets of variables, such that $X, Y, Z \subseteq V(G)$, then X is called **conditionally independent** of Y **given** Z , denoted as

$$X \perp\!\!\!\perp_P Y | Z$$

if and only if

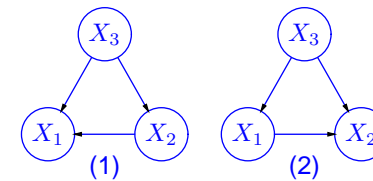
$$P(X | Y, Z) = P(X | Z)$$

Example: Representation of $X_2 \perp\!\!\!\perp_P X_3 | X_1$ in a directed graph



Lecture2: Bayesian networks – p. 23

Chain rule - digraph



Factorisation (1):

$$P(X_1, X_2, X_3) = P(X_1 | X_2, X_3)P(X_2 | X_3)P(X_3)$$

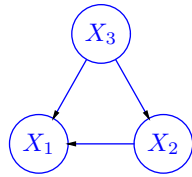
Other factorisation (2):

$$P(X_1, X_2, X_3) = P(X_2 | X_1, X_3)P(X_1 | X_3)P(X_3)$$

\Rightarrow different *factorisations* possible

Lecture2: Bayesian networks – p. 24

Does the chain rule help?



$$P(X_1, X_2, X_3) = P(X_1 | X_2, X_3)P(X_2 | X_3)P(X_3)$$

i.e. we need:

$$\begin{array}{ll}
 P(x_1 | x_2, x_3) & P(x_1 | x_2, \neg x_3) \\
 P(\neg x_1 | x_2, x_3) & P(\neg x_1 | x_2, \neg x_3) \\
 P(x_1 | \neg x_2, x_3) & P(x_1 | \neg x_2, \neg x_3) \\
 P(\neg x_1 | \neg x_2, x_3) & P(\neg x_1 | \neg x_2, \neg x_3) \\
 \vdots & \vdots
 \end{array}$$

Lecture2: Bayesian networks – p. 25

Does the chain rule help?

$$\begin{array}{ll}
 P(x_1, x_2, x_3) & P(x_1, x_2, \neg x_3) \\
 P(\neg x_1, x_2, x_3) & P(\neg x_1, x_2, \neg x_3) \\
 P(x_1, \neg x_2, x_3) & P(x_1, \neg x_2, \neg x_3) \\
 P(\neg x_1, \neg x_2, x_3) & P(\neg x_1, \neg x_2, \neg x_3)
 \end{array}$$

8 required? No, because $\sum_{X_1, X_2, X_3} P(X_1, X_2, X_3) = 1$
Hence, e.g.

$$\begin{aligned}
 P(x_1, x_2, x_3) &= 1 - \sum_{X_2, X_3} P(\neg x_1, X_2, X_3) \\
 &\quad - \sum_{X_3} P(x_1, \neg x_2, X_3) - P(x_1, x_2, \neg x_3)
 \end{aligned}$$

Lecture2: Bayesian networks – p. 27

Does the chain rule help?

$$\begin{array}{ll}
 \vdots & \vdots \\
 P(x_2 | x_3) & P(x_3) \\
 P(\neg x_2 | x_3) & P(\neg x_3) \\
 P(x_2 | \neg x_3) & \\
 P(\neg x_2 | \neg x_3) &
 \end{array}$$

So, 14 probabilities; however

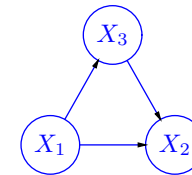
$$\begin{aligned}
 P(x_1 | X_2, X_3) &= 1 - P(\neg x_1 | X_2, X_3), \\
 P(x_2 | X_3) &= 1 - P(\neg x_2 | X_3), \text{ and } P(x_3) = 1 - P(\neg x_3)
 \end{aligned}$$

\Rightarrow 7 probabilities required

How many did we have originally for $P(X_1, X_2, X_3)$?

Lecture2: Bayesian networks – p. 26

Let's use stochastic independence



$$P(X_1, X_2, X_3) = P(X_2 | X_1, X_3)P(X_3 | X_1)P(X_1)$$

Now assume that X_2 and X_3 are **conditionally independent** given X_1 :

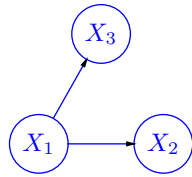
$$P(X_2 | X_1, X_3) = P(X_2 | X_1)$$

and

$$P(X_3 | X_1, X_2) = P(X_3 | X_1)$$

Lecture2: Bayesian networks – p. 28

Stochastic independence: does it help?



$$P(X_2 | X_1, X_3) = P(X_2 | X_1)$$

$$\begin{aligned} P(X_1, X_2, X_3) &= P(X_2 | X_1, X_3)P(X_3 | X_1)P(X_1) \\ &= P(X_2 | X_1)P(X_3 | X_1)P(X_1) \end{aligned}$$

Only $5 = 2 + 2 + 1$ probabilities required instead of 7

Probabilistic inference

Given:

$P(x_4 | x_3) = 0.4$
 $P(x_4 | \neg x_3) = 0.1$
 $P(x_3 | x_1, x_2) = 0.3$
 $P(x_3 | \neg x_1, x_2) = 0.5$
 $P(x_3 | x_1, \neg x_2) = 0.7$
 $P(x_3 | \neg x_1, \neg x_2) = 0.9$
 $P(x_1) = 0.6$
 $P(x_2) = 0.2$

Then:

$$\begin{aligned} P(x_4) &= P(x_4, x_3) + P(x_4, \neg x_3) \\ &\text{(marginalisation)} \\ &= P(x_4 | x_3)P(x_3) + P(x_4 | \neg x_3)P(\neg x_3) \\ &\text{(conditioning)} \\ &= \sum_{X_3} P(x_4 | X_3)P(X_3) \end{aligned}$$

Probabilistic inference

$P(x_4 | x_3) = 0.4$
 $P(x_4 | \neg x_3) = 0.1$
 $P(x_3 | x_1, x_2) = 0.3$
 $P(x_3 | \neg x_1, x_2) = 0.5$
 $P(x_3 | x_1, \neg x_2) = 0.7$
 $P(x_3 | \neg x_1, \neg x_2) = 0.9$
 $P(x_1) = 0.6$
 $P(x_2) = 0.2$

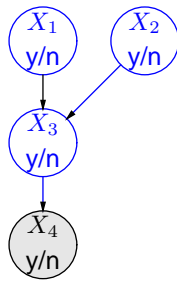
$$\begin{aligned} P(x_3) &= \sum_{X_1, X_2} P(x_3, X_1, X_2) \\ &= \sum_{X_1, X_2} P(x_3 | X_1, X_2)P(X_1, X_2) \\ &= \sum_{X_1, X_2} P(x_3 | X_1, X_2)P(X_1)P(X_2) \\ \Rightarrow P(x_4) &= \sum_{X_3} P(x_4 | X_3) \sum_{X_1, X_2} P(X_3 | X_1, X_2)P(X_1, X_2) \\ &= 0.31 \end{aligned}$$

Probabilistic inference: evidence

$P(x_4 | x_3) = 0.4$
 $P(x_4 | \neg x_3) = 0.1$
 $P(x_3 | x_1, x_2) = 0.3$
 $P(x_3 | \neg x_1, x_2) = 0.5$
 $P(x_3 | x_1, \neg x_2) = 0.7$
 $P(x_3 | \neg x_1, \neg x_2) = 0.9$
 $P(x_1) = 0.6$
 $P(x_2) = 0.2$

$$\begin{aligned} P^*(x_4) = P(x_4 | x_2) &= \sum_{X_3} P(x_4 | x_2, X_3)P(X_3 | x_2) \\ &= \sum_{X_3} P(x_4 | X_3)P(X_3 | x_2) \\ &= \sum_{X_3} P(x_4 | X_3) \sum_{X_1} P(X_3 | X_1, x_2)P(X_1 | x_2) \\ &= \sum_{X_3} P(x_4 | X_3) \sum_{X_1} P(X_3 | X_1, x_2)P(X_1) \\ &= 0.214 \end{aligned}$$

Probabilistic inference: evidence



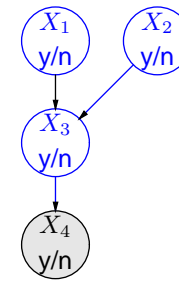
$$\begin{aligned}
 P(x_4 | x_3) &= 0.4 \\
 P(x_4 | \neg x_3) &= 0.1 \\
 P(x_3 | x_1, x_2) &= 0.3 \\
 P(x_3 | \neg x_1, x_2) &= 0.5 \\
 P(x_3 | x_1, \neg x_2) &= 0.7 \\
 P(x_3 | \neg x_1, \neg x_2) &= 0.9 \\
 P(x_1) &= 0.6 \\
 P(x_2) &= 0.2
 \end{aligned}$$

$$P^*(x_2) = P(x_2 | x_4) = \frac{P(x_4 | x_2)P(x_2)}{P(x_4)} \text{ (Bayes' rule)}$$

$$\begin{aligned}
 P(x_4 | x_2) &= 0.214, & P(x_4) &= 0.31 \\
 \Rightarrow P^*(x_2) &= 0.214 \cdot 0.2 / 0.31 \\
 &\approx 0.1381
 \end{aligned}$$

Lecture2: Bayesian networks – p. 33

Probabilistic inference: evidence



$$\begin{aligned}
 P(x_4 | x_3) &= 0.4 \\
 P(x_4 | \neg x_3) &= 0.1 \\
 P(x_3 | x_1, x_2) &= 0.3 \\
 P(x_3 | \neg x_1, x_2) &= 0.5 \\
 P(x_3 | x_1, \neg x_2) &= 0.7 \\
 P(x_3 | \neg x_1, \neg x_2) &= 0.9 \\
 P(x_1) &= 0.6 \\
 P(x_2) &= 0.2
 \end{aligned}$$

$$\begin{aligned}
 P^*(x_2) &= P(x_2 | x_4) = P(x_4 | x_2)P(x_2) / P(x_4) \\
 &= \frac{\sum_{X_3} P(x_4 | X_3) \sum_{X_1} P(X_3 | X_1, x_2) P(X_1) P(x_2)}{\sum_{X_3} P(x_4 | X_3) \sum_{X_1, X_2} P(X_3 | X_1, X_2) P(X_1) P(X_2)}
 \end{aligned}$$

Lecture2: Bayesian networks – p. 34

Bayesian software and links

- Some **software companies** in this area:
 - Hugin (Denmark): www.hugin.dk
 - Norsys (USA): www.norsys.com
 - Knowledge Industries (USA): www.kic.com
 - BayesWare (USA): www.bayesware.com
 - Bayesia (France): www.bayesia.com
- Some **public domain software**:
 - JavaBayes: www-2.cs.cmu.edu/~javabayes
 - BayesBuilder: www.snn.ru.nl/nijmegen/
 - Ideal: Lisp based
 - Elvira: leo.ugr.es/~elvira
 - BNT (Murphy's toolbox, Matlab based): <http://bnt.sourceforge.net/>

Lecture2: Bayesian networks – p. 35