Probabilistic Reasoning Uncertainty and Bayesian networks

Peter Lucas

Institute for Computing and Information Sciences Radboud University Nijmegen, The Netherlands



Uncertainty in Daily Life

• Empirical evidence:

"If symptoms of fever, shortness of breath (dyspnoea), and coughing are present, and the patient has recently visited China, then the patient has *probably* SARS"

• Subjective belief:

ALWEER SARS SLACHTOFFER

"The Balkenende IV government is *likely* to resign soon"

• Temporal dimension:

"There is more than *10% chance* that the Dutch economy will collaps in the next two years"

Uncertainty Representation and Manipulation

- Methods for dealing with uncertainty are **not** new:
 - 17th century: Fermat, Pascal, Huygens, Leibniz, Bernoulli
 - 18th century: Laplace, De Moivre, Bayes
 - 19th century: Gauss, Boole
 - \Rightarrow you could have contributed too if you had been around
- Most important research question in AI:
 - 1970–1987: How to incorporate uncertainty reasoning into logical deduction?
 - 2000-present: How to incorporate uncertainty into logical deduction and induction?

Early AI Methods of Uncertainty

Rule-based uncertainty representation:

fever \land *dyspnoea* \Rightarrow SARS_{CF=0.4}

- Uncertainty calculus (certainty-factor (CF) model, subjective Bayesian method):
 - CF(*fever*, B) = 0.6; CF(*dyspnoea*, B) = 1
 (B is background knowledge)

- Combination functions: $CF(SARS, \{fever, dyspnoea\} \cup B)$ $= 0.4 \cdot max\{0, min\{CF(fever, B), CF(dyspnoea, B)\}\}$ $= 0.4 \cdot max\{0, min\{0.6, 1\}\} = 0.24$ However ···

fever \land *dyspnoea* \Rightarrow SARS_{CF=0.4}

- How likely is the occurrence of *fever* or *dyspnoea* given that the patient has *SARS*?
- How likely is the occurrence of *fever* or *dyspnoea* in the absence of *SARS*?
- How likely is the presence of SARS when just fever is present?
- How likely is *no SARS* when just *fever* is present?







Bayesian Network Formally

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

- G = (V(G), A(G)) is an acyclic directed graph, with
 - $-V(G) = \{X_1, X_2, \dots, X_n\}, \text{ a set of vertices (nodes)}; X \in V(G) \text{ corresponds to a random variable } X$
 - $-A(G) \subseteq V(G) \times V(G)$ a set of arcs reflecting (conditional) independences among variables
- Pr : $\wp(V(G)) \rightarrow [0,1]$ is a joint probability distribution, such that

$$\Pr(V(G)) = \prod_{i=1}^{n} \Pr(X_i \mid \pi_G(X_i))$$

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

Factorisation



Conditional probability distribution:

$$\Pr(X_1 \mid X_2, X_3) = \frac{\Pr(X_1, X_2, X_3)}{\Pr(X_2, X_3)}$$

 $\Rightarrow \Pr(X_1, X_2, X_3) = \Pr(X_1 \mid X_2, X_3) \Pr(X_2 \mid X_3) \Pr(X_3)$

Chain rule yields a factorisation:

$$\Pr(\bigwedge_{i=1}^{n} X_i) = \prod_{i=1}^{n} \Pr(X_i \mid \bigwedge_{k=i+1}^{n} X_k)$$

Independence Representation in Graphs

The set of variables X is conditionally independent of the set Z given the set Y, notation $X \perp\!\!\!\perp Z \mid Y$, iff

 $\Pr(X \mid Y, Z) = \Pr(X \mid Y)$

Three flavours of graph-representation of (in)dependence:

Diverging: Y blocks X and Z: $X \perp \!\!\!\perp Z \mid Y$



Serial: Y blocks X and Z: $X \perp \!\!\!\perp Z \mid Y$



Converging: Y connects X and Z: $X \not\perp Z \mid Y$



Use of Independence Information

General:

 $\Pr(X_1, X_2, X_3) = \Pr(X_2 \mid X_1, X_3) \Pr(X_3 \mid X_1) \Pr(X_1)$

Assume that $X_2 \perp X_3 \mid X_1$, then: $\Pr(X_2 \mid X_1, X_3) = \Pr(X_2 \mid X_1)$ and

 $\Pr(X_3 \mid X_1, X_2) = \Pr(X_3 \mid X_1)$



Only 5 = 2 + 2 + 1 probabilities needed for $Pr(X_1, X_2, X_3)$ (instead of 7)



- SARS ⊥ TEMP | FEVER
- VisitToChina ⊥⊥ DYSPNOEA | SARS

Probabilistic Reasoning

• Interested in marginal probability distributions:

 $\Pr(V_i \mid \mathcal{E}) = \Pr^{\mathcal{E}}(V_i)$

for (possibly empty) evidence \mathcal{E} (instantiated variables)

- Joint probability distribution Pr(V):
 - marginalisation:

$$\Pr(\mathcal{W}) = \sum_{V \setminus \mathcal{W}} \Pr(V)$$
$$= \sum_{V \setminus \mathcal{W}} \prod_{X \in V} \Pr(X \mid \pi(X))$$

- conditional probabilities and Bayes' rule:

$$\Pr(Y, Z \mid X) = \frac{\Pr(X \mid Y, Z) \Pr(Y, Z)}{\Pr(X)}$$

• Many efficient Bayesian reasoning algorithms exist

Naive Probabilistic Reasoning: Evidence



 $\Pr^{\mathcal{E}}(x_{2}) = \Pr(x_{2} \mid x_{4}) = \frac{\Pr(x_{4} \mid x_{2}) \Pr(x_{2})}{\Pr(x_{4})} \text{ (Bayes' rule)}$ $= \frac{\sum_{X_{3}} \Pr(x_{4} \mid X_{3}) \sum_{X_{1}} \Pr(X_{3} \mid X_{1}, x_{2}) \Pr(X_{1}) \Pr(x_{2})}{\sum_{X_{3}} \Pr(x_{4} \mid X_{3}) \sum_{X_{1}, X_{2}} \Pr(X_{3} \mid X_{1}, X_{2}) \Pr(X_{1}) \Pr(X_{2})}$ ≈ 0.14



- Object-oriented approach: vertices are objects, which have local information and carry out local computations
- Updating of probability distribution by message passing: arcs are communication channels



- $\mathcal{E} = \mathcal{E}_{V_i}^+ \cup \mathcal{E}_{V_i}^-$: evidence
- α : normalisation constant

Entering Observations

Observed joint probability distribution:

 $\mathsf{Pr}^{\mathcal{E}}(V) := \mathsf{Pr}(V \mid \mathcal{E})$

- \mathcal{E} : observed random variables
- $U = V \setminus \mathcal{E}$: unobserved random variables

Graphical consequences of observations:

- Additional (observed) dependences: moral lines
- Additional (observed) independences: observed and semiobserved arcs

Observation Transformation

- Moral lines: connect non-connected parents of an observed (descendant of) a common child
- Arc-line transformation: change *arcs* of ancestors of the observed variables into *lines*
- Delete (semi)observed arcs
- Example:



Problem Solving

Bayesian networks are declarative, i.e.:

- mathematical basis
- problem to be solved determined by (1) entered evidence
 E (may include decisions); (2) given hypothesis *H*:

 $\mathsf{Pr}(H \mid \mathcal{E}) \qquad \qquad (\mathsf{cf. KB} \land H \vDash \mathcal{E})$

Examples:

- Description of populations
- Classification and diagnosis: $D = \arg \max_{H} \Pr(H | \mathcal{E})$
- Temporal reasoning, prediction, what-if scenarios
- Decision-making based on decision theory

$$\mathsf{MEU}(D \mid \mathcal{E}) = \max_{d \in D} \sum_{x \in X_{\pi(U)}} u(x) \operatorname{Pr}(x \mid d, \mathcal{E})$$



People become colonised by bacteria when entering a hospital, which may give rise to infection

Bayesian-network Modelling

Qualitative

Quantitative

causal modelling interaction modelling

 $Cause \rightarrow Effect$









e if $P(e | I_1, ..., I_n) = 1$



- Interactions among causes: logical OR
- Meaning: presence of any one of the causes C_i with absolute certainty will cause the effect e (i.e. E = true)

$$\Pr(e|C_1, C_2) = \sum_{\substack{I_1 \lor I_2 = e \\ Pr(i_1|C_1) \Pr(i_2|C_2)}} \Pr(I_k | C_k)$$

=
$$\Pr(i_1|C_1) \Pr(i_2|C_2) + \Pr(\neg i_1|C_1) \Pr(i_2|C_2)$$

+
$$\Pr(i_1|C_1) \Pr(\neg i_2|C_2)$$

• Assessment of O(n) instead of $O(2^n)$ probabilities

Example BN: non-Hodgkin Lymphoma



Bayesian Network Learning

Bayesian network $\mathcal{B} = (G, \mathsf{Pr})$, with

- digraph G = (V(G), A(G)), and
- probability distribution Pr



Learning Bayesian Networks

Problems:

- for many BNs too many probabilities have to be assessed
- complex BNs do not necessarily yield better classifiers
- complex BNs may yield better estimates of a probability distribution

Solution:

- use simple probabilistic models for classification:
 - naive (independent) form BN
 - Tree-Augmented Bayesian Network (TAN)
 - Forest-Augmented Bayesian Network (FAN)
- use background knowledge and clever heuristics

Naive (independent) form BN



- C is a class variable
- The evidence variables E_i in the evidence $\mathcal{E} \subseteq \{E_1, \ldots, E_m\}$ are conditionally independent given the class variable C

This yields:

$$P(C \mid \mathcal{E}) = \frac{P(\mathcal{E} \mid C)P(C)}{P(\mathcal{E})} = \frac{\prod_{E \in \mathcal{E}} P(E \mid C)}{\sum_{C} P(\mathcal{E} \mid C)P(C)}$$

as $E_i \perp E_j \mid C$, for $i \neq j$

Classifier: $c_{\max} = \arg \max_C P(C \mid \mathcal{E})$

Learning Structure from Data

Given the following dataset D:

Student	Gender	IQ	High Mark for Maths
1	male	low	no
2	female	average	yes
3	male	high	yes
4	female	high	yes

and the following Bayesian networks:







 $Q(G,D) = \log \Pr(G) - |D| \cdot H(G,D) - \frac{1}{2}k \cdot \log |D|$, where:

- Pr(G): prior probability of G
- -H(G,D): negative value of match
- $-\frac{1}{2}k \cdot \log |D|$: penalty term

Research Issues

Qualitative modelling:

- To determine the structure of a network
- Assessment of $\Pr(V_i \mid \pi(V_i))$

Probabilistic-logic learning

 \mathbf{BR}_{C}

 (BR_B)

Inf

BR

- Structure learning: determine the 'best' graph topology
- Parameter learning: determine the 'best' probability distribution (discrete or continuous)
- Bayesian (probabilistic) logic and relational learning

 \Rightarrow you can contribute too \cdots