

Week 1: Quantifying Uncertainty & Bayes Network

1. Quantifying Uncertainty

1.1. Introduction

- 1.1.1. Definition of Agent
- 1.1.2. Types of Uncertainty
- 1.1.3. Why cares about uncertainty
- 1.1.4. Traditional way
- 1.1.5. Utility Theory & Decision Theory
- 1.1.5. Why probability:
- 1.1.6. Basic Probability:
 - 1.1.6.1. Full Joint Distributions:
 - 1.1.6.2. Absolute independence:
 - 1.1.6.3. Conditional Independence :
 - 1.1.6.4. Bayes Rule :

2. Probabilistic Reasoning

2.1. Bayesian network

- 2.1.1. Rules:
- 2.1.2. Semantics
 - 2.1.2.1. Representing full joint distribution
 - 2.1.2.2. conditional independence relations
- 2.1.3. Ways of constructing

2.2. Exact Inference in Bayesian Networks

- 2.2.1. Variables
- 2.2.2. Inference by enumeration
- 2.2.3. Variable Elimination Algorithm

2.3. Approximate Inference

- 2.3.1. Ancestral Sampling
- 2.3.2. Rejection Sampling
- 2.3.3. Likelihood Weighting

2.4. Markov chain Monte Carlo (MCMC)

1. Quantifying Uncertainty

1.1. Introduction

1.1.1. Definition of Agent

1. sensors → Perceiving Environment
2. Actuators → actions

1.1.2. Types of Uncertainty

1. stochasticity

- a. outcome uncertainty
- b. environment fluctuation
- 2. state uncertainty
 - a. sensor limited
 - b. noise

1.1.3. Why cares about uncertainty

an agent may never know **its state/sequence of actions** exactly

1.1.4. Traditional way

Belief state → set of all possible states (**no probability**)

1. cons: **qualification problems**: needs to specify all exceptions
2. solution: **rational decision**: goals+likelihood+achieved degree (why we need **probability** tools)

1.1.5. Utility Theory & Decision Theory

1. **Utility Theory** → preference
2. Decision Theory = probability theory + utility theory
3. **rational** agent ⇔ it chooses the highest expected utility (**MEU**)
4. Expected Utility:
 - a. outcomes uncertainty:

$$U(a) = \sum_o u(o) * b(o|a)$$

$$b. \text{ states uncertainty: } U(a) = \sum_s u(s, a) * b(s)$$

1.1.5. Why probability:

But de Finetti Argument: “If Agent 1 expresses a set of degrees of belief that violate the axioms of probability theory then there is a combination of bets by Agent 2 that guarantees that Agent 1 will lose money every time”

1.1.6. Basic Probability:

1.1.6.1. Full Joint Distributions:

cons: Too large

1.1.6.2. Absolute independence:

pros: helpful when abstract distributions $2^n \rightarrow n$

cons: hard to meet

1.1.6.3. Conditional Independence :

$$\begin{aligned}
P(X, Y|Z) &= P(X|Z)P(Y|Z) \\
\text{or } \mathbf{P}(X | Y, Z) &= \mathbf{P}(X | Z) \\
\text{or } \mathbf{P}(Y | X, Z) &= \mathbf{P}(Y | Z)
\end{aligned}$$

1.1.6.4. Bayes Rule :

$$P(b|a) = \frac{P(a|b)P(b)}{P(a)} = \frac{P(a|b)P(b)}{\sum P(a|b_i)p(b_i)}$$

$$P(Y|X, e) = P(X|Y, e)P(Y|e)/P(X|e)$$

$$P(Y|X) = \alpha P(X|Y)P(Y)$$

1. enable update probability
2. “diagnostic knowledge is often more fragile than causal knowledge”

2. Probabilistic Reasoning

2.1. Bayesian network

Bayesian network (compact representation) \Leftrightarrow any joint probability distribution

2.1.1. Rules:

1. node \rightarrow random variable
2. If there is an arrow from node X to node Y, X is said to be a **parent** of Y, intuitively meaning X has a **directly influence** on Y
3. Each node has a conditional probability distribution $P(X_i | Parents(X_i))$
4. Acyclic

2.1.2. Semantics

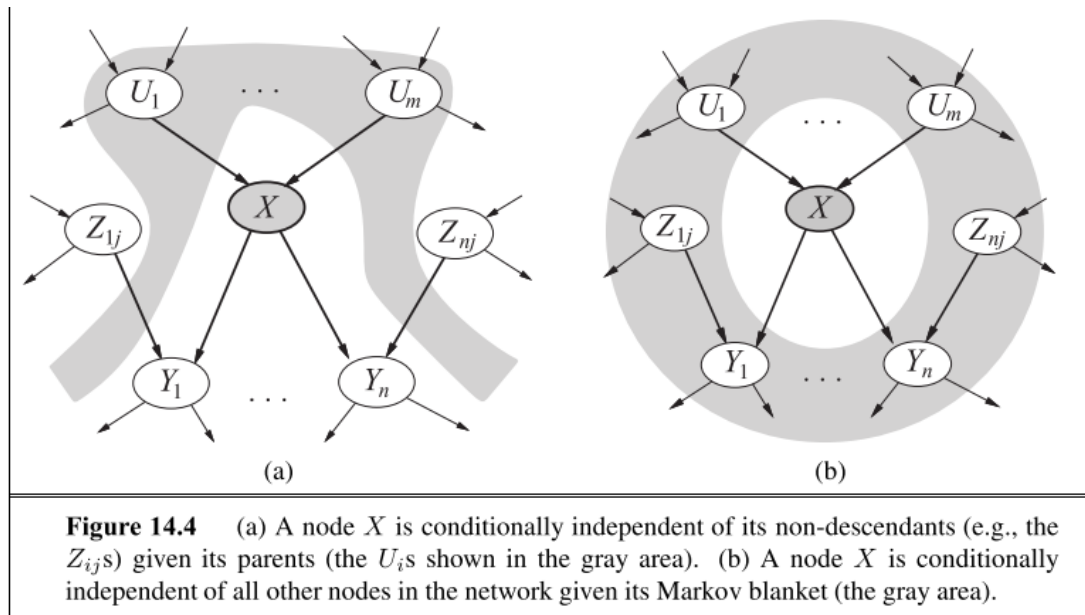
2.1.2.1. Representing full joint distribution

the conditional probability distribution ensured that the Bayesian Network can representing full joint distribution comprehensively

$$P(x_1, x_2, \dots, x_n) = \prod_{i=1}^n P(x_i | Parents(X_i))$$

2.1.2.2. conditional independence relations

1. given its **parents**, each variables is conditionally independent of its **non-descendants**
2. given its **Markov blanket**(parents, children, children’s parents) \rightarrow all other nodes



2.1.3. Ways of constructing

Based on chain rule:

$$\mathbf{P}(X_i | X_{i-1}, \dots, X_1) = \mathbf{P}(X_i | \text{Parents}(X_i)) \quad (14.3)$$

1. *Nodes*: First determine the set of variables that are required to model the domain. Now order them, $\{X_1, \dots, X_n\}$. Any order will work, but the resulting network will be more compact if the variables are ordered such that causes precede effects.
2. *Links*: For $i = 1$ to n do:
 - Choose, from X_1, \dots, X_{i-1} , a minimal set of parents for X_i , such that Equation (14.3) is satisfied.
 - For each parent insert a link from the parent to X_i .
 - CPTs: Write down the conditional probability table, $\mathbf{P}(X_i | \text{Parents}(X_i))$.

2.2. Exact Inference in Bayesian Networks

2.2.1. Variables

$$P(X|e)$$

1. Query Variables: X ;
2. Evidence Variables: Always the condition variables, e ;
3. Hidden Variables: non-evidence non query variables;

2.2.2. Inference by enumeration

1. Using Law of total cumulance:

$$P(X|e) = \alpha P(X, e) = \alpha \sum_y P(X, e, y)$$

2. Using Bayes Network:

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

$$\Rightarrow \text{using } P(x_1, x_2, \dots, x_n) = \prod_{i=1}^n P(x_i | \text{Parents}(X_i))$$

2.2.3. Variable Elimination Algorithm

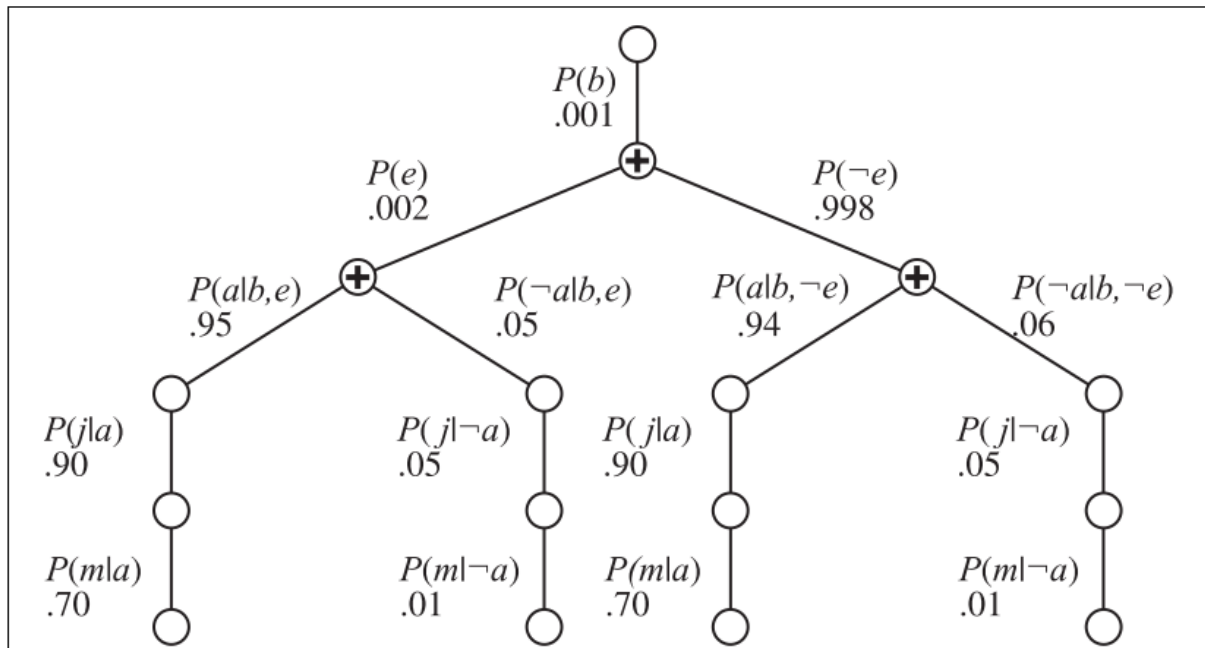


Figure 14.8 The structure of the expression shown in Equation (14.4). The evaluation proceeds top down, multiplying values along each path and summing at the “+” nodes. Notice the repetition of the paths for j and m .

By using pointwise product of a pair of factors, and summing out a variable from a product of factors:

- First push in the summations as far as possible
- Then carry out summations **right-to-left**, caching intermediate results in new **factors**

$$\begin{aligned}
\mathbf{P}(B|j, m) &= \alpha \underbrace{\mathbf{P}(B)}_B \underbrace{\sum_e P(e)}_E \underbrace{\sum_a \mathbf{P}(a|B, e)}_A \underbrace{P(j|a)}_J \underbrace{P(m|a)}_M \\
&= \alpha \mathbf{P}(B) \sum_e P(e) \sum_a \mathbf{P}(a|B, e) P(j|a) f_M(a) \\
&= \alpha \mathbf{P}(B) \sum_e P(e) \sum_a \mathbf{P}(a|B, e) f_J(a) f_M(a) \\
&= \alpha \mathbf{P}(B) \sum_e P(e) \sum_a f_A(a, b, e) f_J(a) f_M(a) \\
&= \alpha \mathbf{P}(B) \sum_e P(e) f_{\bar{A}JM}(b, e) \text{ (sum out } A) \\
&= \alpha \mathbf{P}(B) f_{\bar{E}\bar{A}JM}(b) \text{ (sum out } E) \\
&= \alpha f_B(b) \times f_{\bar{E}\bar{A}JM}(b)
\end{aligned}$$

The time and space requirements of variable elimination are dominated by the size of the largest factor constructed during the operation of the algorithm

One fairly effective method is a **greedy one**: eliminate whichever variable minimizes the size of the next factor to be constructed

2.3. Approximate Inference

Typically two ways:

1. Based on sampling
2. Based on optimization

We only covered method based on sampling

2.3.1. Ancestral Sampling

The idea is to sample each variable in turn, in topological order

```

function PRIOR-SAMPLE(bn) returns an event sampled from bn
  inputs: bn, a belief network specifying joint distribution  $P(X_1, \dots, X_n)$ 

   $\mathbf{x} \leftarrow$  an event with  $n$  elements
  for  $i = 1$  to  $n$  do
     $x_i \leftarrow$  a random sample from  $P(X_i \mid \text{parents}(X_i))$ 
      given the values of  $\text{Parents}(X_i)$  in  $\mathbf{x}$ 
  return  $\mathbf{x}$ 

```

The consistency can be proved from law of large numbers:

Probability that PRIORSAMPLE generates a particular event

$$S_{PS}(x_1 \dots x_n) = \prod_{i=1}^n P(x_i \mid \text{parents}(X_i)) = P(x_1 \dots x_n)$$

i.e., the true prior probability

E.g., $S_{PS}(t, f, t, t) = 0.5 \times 0.9 \times 0.8 \times 0.9 = 0.324 = P(t, f, t, t)$

Let $N_{PS}(x_1 \dots x_n)$ be the number of samples generated for event x_1, \dots, x_n

Then we have

$$\begin{aligned}
 \lim_{N \rightarrow \infty} \hat{P}(x_1, \dots, x_n) &= \lim_{N \rightarrow \infty} N_{PS}(x_1, \dots, x_n) / N && \{\text{law of large numbers}\} \\
 &= S_{PS}(x_1, \dots, x_n) \\
 &= P(x_1 \dots x_n)
 \end{aligned}$$

That is, estimates derived from PRIORSAMPLE are consistent

Shorthand: $\hat{P}(x_1, \dots, x_n) \approx P(x_1 \dots x_n)$

2.3.2. Rejection Sampling

Ancestral Sampling do not care about observations variables. Rejection Sampling Method is to estimate something more complicated... like $P(X|e)$.

```

function REJECTION-SAMPLING( $X, e, bn, N$ ) returns an estimate of  $P(X|e)$ 
  local variables:  $N$ , a vector of counts over  $X$ , initially zero
  for  $j = 1$  to  $N$  do
     $x \leftarrow \text{PRIOR-SAMPLE}(bn)$ 
    if  $x$  is consistent with  $e$  then
       $N[x] \leftarrow N[x] + 1$  where  $x$  is the value of  $X$  in  $x$ 
  return NORMALIZE( $N[X]$ )

```

eg. $\{ \begin{array}{l} N[x_1=1, x_2=1] \\ N[x_1=0, x_2=1] \end{array} \}$

E.g., estimate $P(\text{Rain}|\text{Sprinkler}=\text{true})$ using 100 samples

27 samples have $\text{Sprinkler}=\text{true}$

Of these, 8 have $\text{Rain}=\text{true}$ and 19 have $\text{Rain}=\text{false}$.

$$\hat{P}(\text{Rain}|\text{Sprinkler}=\text{true}) = \text{NORMALIZE}(\langle 8, 19 \rangle) = \langle 0.296, 0.704 \rangle$$

The consistency can be proved from law of large numbers:

$$\begin{aligned}
 \hat{P}(X|e) &= \alpha N_{PS}(X, e) && \text{(algorithm defn.)} \\
 &= N_{PS}(X, e) / N_{PS}(e) && \text{(normalized by } N_{PS}(e)) \\
 &\approx P(X, e) / P(e) && \text{(property of PRIOR-SAMPLE)} \\
 &= P(X|e) && \text{(defn. of conditional probability)}
 \end{aligned}$$

Hence rejection sampling **returns consistent posterior estimates**

Cons: Hopeless Expensive if $P(e)$ is small, and $P(e)$ does drop off exponentially with number of evidence variables.

2.3.3. Likelihood Weighting

Likelihood Weighting Method only sample points consistent with evidence e by using weight.


```

function LIKELIHOOD-WEIGHTING( $X, e, bn, N$ ) returns an estimate of  $P(X|e)$ 
  local variables:  $\mathbf{W}$ , a vector of weighted counts over  $X$ , initially zero
  for  $j = 1$  to  $N$  do
     $\mathbf{x}, w \leftarrow \text{WEIGHTED-SAMPLE}(bn)$ 
     $\mathbf{W}[x] \leftarrow \mathbf{W}[x] + w$  where  $x$  is the value of  $X$  in  $\mathbf{x}$ 
  return NORMALIZE( $\mathbf{W}[X]$ )

```

```

function WEIGHTED-SAMPLE( $bn, e$ ) returns an event and a weight
   $\mathbf{x} \leftarrow$  an event with  $n$  elements;  $w \leftarrow 1$ 
  for  $i = 1$  to  $n$  do
    if  $X_i$  has a value  $x_i$  in  $e$ 
      then  $w \leftarrow w \times P(X_i = x_i \mid \text{parents}(X_i))$ 
      else  $x_i \leftarrow$  a random sample from  $P(X_i \mid \text{parents}(X_i))$ 
  return  $\mathbf{x}, w$ 

```

The consistency can be proved from law of large numbers:

■ Can again show this is consistent:

Sampling probability for WEIGHTEDSAMPLE is

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

Note: pays attention to evidence in **ancestors** only

\Rightarrow somewhere “in between” prior and posterior distribution

Weight for a given sample \mathbf{z}, \mathbf{e} is

$$w(\mathbf{z}, \mathbf{e}) = \prod_{i=1}^m P(e_i | \text{parents}(E_i))$$

Weighted sampling probability is

$$\begin{aligned}
 & S_{WS}(\mathbf{z}, \mathbf{e}) w(\mathbf{z}, \mathbf{e}) \\
 &= \prod_{i=1}^l P(z_i | \text{parents}(Z_i)) \prod_{i=1}^m P(e_i | \text{parents}(E_i)) \\
 &= P(\mathbf{z}, \mathbf{e}) \text{ (by standard global semantics of network)}
 \end{aligned}$$



cons:

1. So, this algorithms may suffer degradation in performance as the number of evidence variables increases. This is because **most samples will have very low weights** and hence the weighted estimate will **be dominated by the tiny fraction of samples** that accord more than an infinitesimal likelihood to the evidence

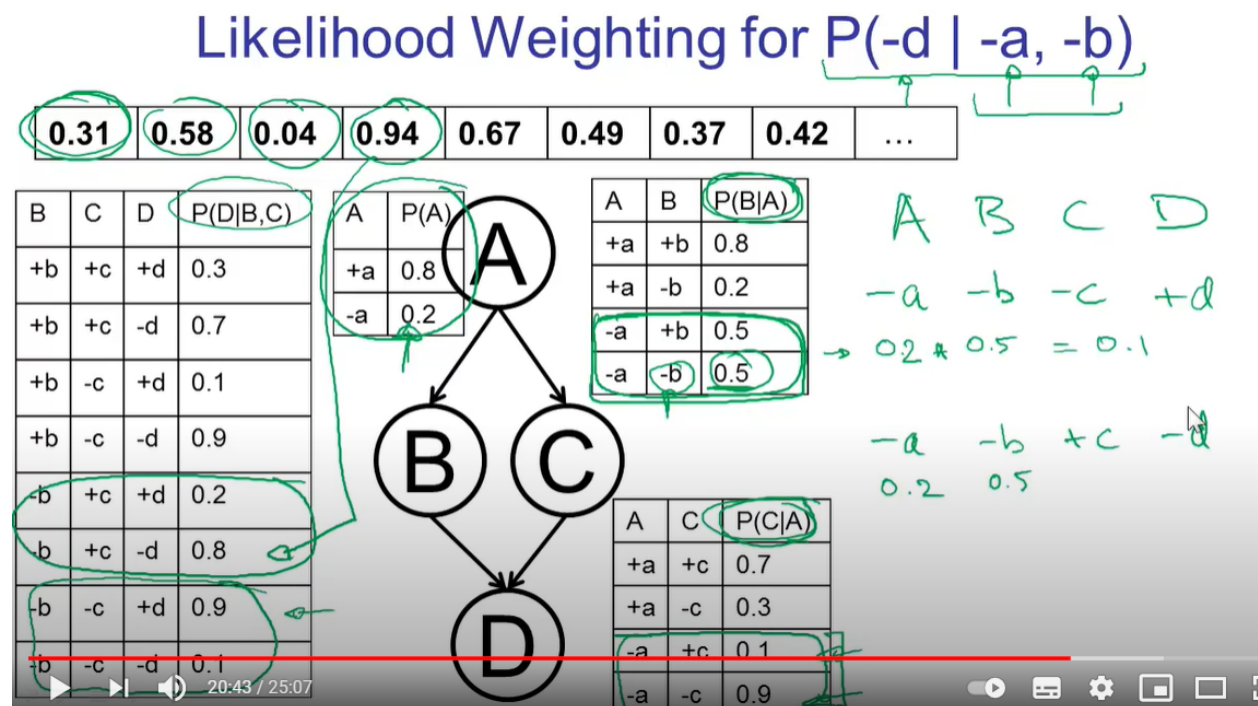
2. The evidence based on which we sample **influences the downstream random variables** that we sample and not upstream. This could bias the samples that are generated.
3. If the evidence variables occur late in the variable ordering, the nonevidence variables will have no evidence in their parents and ancestors to guide the generation of samples.

Intuition:

Based on the conditionally independent in the Bayesian Network and limitedly taken known information (observation) into consideration

Example

https://www.youtube.com/watch?v=ol0l6aTfb_g



2.4. Markov chain Monte Carlo (MCMC)

Waiting for updating