

Inference in belief networks

CHAPTER 15.3-4 + NEW

Inference tasks

- Simple queries:** compute posterior marginal $\mathbf{P}(X_i|\mathbf{E}=\mathbf{e})$
e.g., $P(\text{NoGas}|\text{Gauge} = \text{empty}, \text{Lights} = \text{on}, \text{Starts} = \text{false})$
- Conjunctive queries:** $\mathbf{P}(X_i, X_j|\mathbf{E}=\mathbf{e}) = \mathbf{P}(X_i|\mathbf{E}=\mathbf{e})\mathbf{P}(X_j|X_i, \mathbf{E}=\mathbf{e})$
- Optimal decisions:** decision networks include utility information;
probabilistic inference required for $P(\text{outcome}|\text{action}, \text{evidence})$
- Value of information:** which evidence to seek next?
- Sensitivity analysis:** which probability values are most critical?
- Explanation:** why do I need a new starter motor?

Enumeration algorithm

Exhaustive depth-first enumeration: $O(n)$ space, $O(d^n)$ time

```

ENUMERATION ASK( $X, e, bn$ ) returns a distribution over  $X$ 
inputs:  $X$ , the query variable
        $e$ , evidence specified as an event
        $bn$ , a belief network specifying joint distribution  $\mathbf{P}(X_1, \dots, X_n)$ 

 $Q(x) \leftarrow$  a distribution over  $X$ 
for each value  $x_i$  of  $X$  do
    extend  $e$  with value  $x_i$  for  $X$ 
     $Q(x_i) \leftarrow$  ENUMERATE ALL(VARS[bn],  $e$ )
return NORMALIZE( $Q(X)$ )

ENUMERATE ALL( $vars, e$ ) returns a real number
if EMPTY?( $vars$ ) then return 1.0
else do
     $Y \leftarrow$  FIRST( $vars$ )
    if  $Y$  has value  $y$  in  $e$ 
        then return  $P(y | Pa(Y)) \times$  ENUMERATE ALL(REST( $vars$ ),  $e$ )
    else return  $\sum_y P(y | Pa(Y)) \times$  ENUMERATE ALL(REST( $vars$ ),  $e_y$ )
        where  $e_y$  is  $e$  extended with  $Y = y$ 
    
```

Outline

- ◇ Exact inference by enumeration
- ◇ Exact inference by variable elimination
- ◇ Approximate inference by stochastic simulation
- ◇ Approximate inference by Markov chain Monte Carlo

Inference by enumeration

- Slightly intelligent way to sum out variables from the joint without actually constructing its explicit representation
- Simple query on the burglary network:
 $\mathbf{P}(B|J = \text{true}, M = \text{true})$
 $= \mathbf{P}(B, J = \text{true}, M = \text{true}) / \mathbf{P}(J = \text{true}, M = \text{true})$
 $= \alpha \mathbf{P}(B, J = \text{true}, M = \text{true})$
 $= \alpha \sum_e \sum_a \mathbf{P}(B, e, a, J = \text{true}, M = \text{true})$
- Rewrite full joint entries using product of CPT entries:
 $\mathbf{P}(B = \text{true} | J = \text{true}, M = \text{true})$
 $= \alpha \sum_e \sum_a \mathbf{P}(B = \text{true}) \mathbf{P}(e) \mathbf{P}(a | B = \text{true}, e) \mathbf{P}(J = \text{true} | a) \mathbf{P}(M = \text{true} | a)$
 $= \alpha \mathbf{P}(B = \text{true}) \sum_e \mathbf{P}(e) \sum_a \mathbf{P}(a | B = \text{true}, e) \mathbf{P}(J = \text{true} | a) \mathbf{P}(M = \text{true} | a)$

Inference by variable elimination

- Enumeration is inefficient: repeated computation
e.g., computes $P(J = \text{true} | a) P(M = \text{true} | a)$ for each value of e
- Variable elimination: carry out summations right-to-left, storing intermediate results (**factors**) to avoid recomputation
- $\mathbf{P}(B | J = \text{true}, M = \text{true})$
 $= \alpha \underbrace{\mathbf{P}(B)}_B \underbrace{\sum_e \mathbf{P}(e)}_E \underbrace{\sum_a \mathbf{P}(a | B, e)}_A \underbrace{\mathbf{P}(J = \text{true} | a)}_J \underbrace{\mathbf{P}(M = \text{true} | a)}_M$
 $= \alpha \mathbf{P}(B) \sum_e \mathbf{P}(e) \sum_a \mathbf{P}(a | B, e) \mathbf{P}(J = \text{true} | a) f_M(a)$
 $= \alpha \mathbf{P}(B) \sum_e \mathbf{P}(e) \sum_a \mathbf{P}(a | B, e) f_J(a) f_M(a)$
 $= \alpha \mathbf{P}(B) \sum_e \mathbf{P}(e) \sum_a f_{AJM}(a, b, e) f_J(a) f_M(a)$
 $= \alpha \mathbf{P}(B) \sum_e \mathbf{P}(e) f_{AJM}(b, e)$ (sum out A)
 $= \alpha \mathbf{P}(B) f_{EAJM}(b)$ (sum out E)
 $= \alpha f_B(b) \times f_{EAJM}(b)$

Variable elimination: Basic operations

Pointwise product of factors f_1 and f_2 :

$$f_1(x_1, \dots, x_j, y_1, \dots, y_k) \times f_2(y_1, \dots, y_k, z_1, \dots, z_l) \\ = f(x_1, \dots, x_j, y_1, \dots, y_k, z_1, \dots, z_l)$$

E.g., $f_1(a, b) \times f_2(b, c) = f(a, b, c)$

Summing out a variable from a product of factors: move any constant factors outside the summation:

$$\sum_x f_1 \times \dots \times f_k = f_1 \times \dots \times f_i \sum_x f_{i+1} \times \dots \times f_k = f_1 \times \dots \times f_i \times f_X$$

assuming f_1, \dots, f_i do not depend on X

Variable elimination algorithm

function ELIMINATIONASK(X, e, bn) returns a distribution over X

inputs: X , the query variable
 e , evidence specified as an event
 bn , a belief network specifying joint distribution $P(X_1, \dots, X_n)$

if $X \in e$ **then return** observed point distribution for X

$factors \leftarrow []$; $vars \leftarrow REVERSE(VARS[bn])$

for each var **in** $vars$ **do**

$factors \leftarrow [MAKEFACTOR(var, e)]factors$

if var is a hidden variable **then** $factors \leftarrow SUMOUT(var, factors)$

return NORMALIZE(POINTWISEPRODUCT($factors$))

Complexity of exact inference

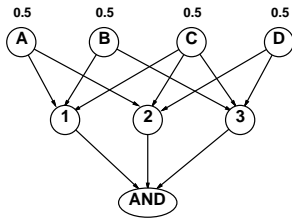
Singly connected networks (or polytrees):

- any two nodes are connected by at most one (undirected) path
- time and space cost of variable elimination are $O(d^k n)$

Multiply connected networks:

- can reduce 3SAT to exact inference \Rightarrow NP-hard
- equivalent to counting 3SAT models \Rightarrow #P-complete

- $A \vee B \vee C$
- $C \vee D \vee \neg A$
- $B \vee C \vee \neg D$



Inference by stochastic simulation

Basic idea:

- Draw N samples from a sampling distribution S
- Compute an approximate posterior probability \hat{P}
- Show this converges to the true probability P

Outline:

- Sampling from an empty network
- Rejection sampling: reject samples disagreeing with evidence
- Likelihood weighting: use evidence to weight samples
- MCMC: sample from a stochastic process whose stationary distribution is the true posterior

Sampling from an empty network

function PRIORSAMPLE(bn) returns an event sampled from $P(X_1, \dots, X_n)$ specified by bn
 $x \leftarrow$ an event with n elements
for $i = 1$ **to** n **do**
 $x_i \leftarrow$ a random sample from $P(X_i | Parents(X_i))$
return x

$P(Cloudy) = \langle 0.5, 0.5 \rangle$

sample \rightarrow true

$P(Sprinkler|Cloudy) = \langle 0.1, 0.9 \rangle$

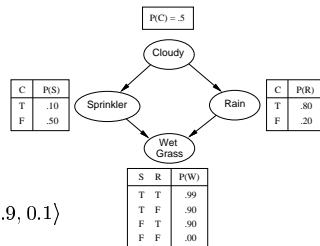
sample \rightarrow false

$P(Rain|Cloudy) = \langle 0.8, 0.2 \rangle$

sample \rightarrow true

$P(WetGrass|\neg Sprinkler, Rain) = \langle 0.9, 0.1 \rangle$

sample \rightarrow true



Sampling from an empty network contd.

Probability that PRIORSAMPLE generates a particular event

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

i.e., the true prior probability

Let $N_{PS}(\mathbf{Y} = \mathbf{y})$ be the number of samples generated for which $\mathbf{Y} = \mathbf{y}$, for any set of variables \mathbf{Y} .

Then $\hat{P}(\mathbf{Y} = \mathbf{y}) = N_{PS}(\mathbf{Y} = \mathbf{y})/N$ and

$$\lim_{N \rightarrow \infty} \hat{P}(\mathbf{Y} = \mathbf{y}) = \sum_{\mathbf{h}} h_{PS}(\mathbf{Y} = \mathbf{y}, \mathbf{H} = \mathbf{h}) \\ = \sum_{\mathbf{h}} h(\mathbf{Y} = \mathbf{y}, \mathbf{H} = \mathbf{h}) \\ = P(\mathbf{Y} = \mathbf{y})$$

That is, estimates derived from PRIORSAMPLE are consistent

Rejection sampling

$\hat{P}(X|e)$ estimated from samples agreeing with e

```

function REJECTION_SAMPLING( $X, e, bn, N$ ) returns an approximation to  $P(X|e)$ 
   $N[X] \leftarrow$  a vector of counts over  $X$ , initially zero
  for  $j = 1$  to  $N$  do
     $x \leftarrow$  PRIORSAMPLE( $bn$ )
    if  $x$  is consistent with  $e$  then
       $N[x] \leftarrow N[x] + 1$  where  $x$  is the value of  $X$  in  $x$ 
  return NORMALIZE( $N[X]$ )
  
```

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 $\text{Rain} = \text{false}$.

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

Similar to a basic real-world empirical estimation procedure

Analysis of rejection sampling

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

Hence rejection sampling returns consistent posterior estimates

Problem: hopelessly expensive if $P(e)$ is small

Likelihood weighting

Idea: fix evidence variables, sample only nonevidence variables, and weight each sample by the likelihood it accords the evidence

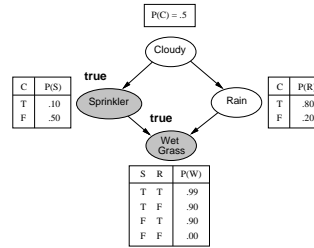
```

function WEIGHTEDSAMPLE( $bn, e$ ) returns an event and a weight
   $x \leftarrow$  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  $w \leftarrow w \times P(X_i = x_i | \text{Parents}(X_i))$ 
      else  $x_i \leftarrow$  a random sample from  $P(X_i | \text{Parents}(X_i))$ 
  return  $x, w$ 

function LIKELIHOODWEIGHTING( $X, e, bn, N$ ) returns an approximation to  $P(X|e)$ 
   $W[X] \leftarrow$  a vector of weighted counts over  $X$ , initially zero
  for  $j = 1$  to  $N$  do
     $x, w \leftarrow$  WEIGHTEDSAMPLE( $bn$ )
     $W[x] \leftarrow W[x] + w$  where  $x$  is the value of  $X$  in  $x$ 
  return NORMALIZE( $W[X]$ )
  
```

Likelihood weighting example

Estimate $P(\text{Rain}|\text{Sprinkler} = \text{true}, \text{WetGrass} = \text{true})$



LW example contd.

Sample generation process:

1. $w \leftarrow 1.0$
2. Sample $P(\text{Cloudy}) = \langle 0.5, 0.5 \rangle$; say *true*
3. *Sprinkler* has value *true*, so
 $w \leftarrow w \times P(\text{Sprinkler} = \text{true} | \text{Cloudy} = \text{true}) = 0.1$
4. Sample $P(\text{Rain} | \text{Cloudy} = \text{true}) = \langle 0.8, 0.2 \rangle$; say *true*
5. *WetGrass* has value *true*, so
 $w \leftarrow w \times P(\text{WetGrass} = \text{true} | \text{Sprinkler} = \text{true}, \text{Rain} = \text{true}) = 0.099$

Likelihood weighting analysis

Sampling probability for WEIGHTEDSAMPLE is

$$S_{WS}(y, e) = \prod_{i=1}^l P(y_i | \text{Parents}(Y_i))$$

Note: pays attention to evidence in *ancestors* only

\Rightarrow somewhere "in-between" prior and posterior distribution

Weight for a given sample y, e is

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

Weighted sampling probability is

$$\begin{aligned} S_{WS}(y, e) w(y, e) &= \prod_{i=1}^l P(y_i | \text{Parents}(Y_i)) \prod_{i=1}^m P(e_i | \text{Parents}(E_i)) \\ &= P(y, e) \text{ (by standard global semantics of network)} \end{aligned}$$

Hence likelihood weighting returns consistent estimates
 but performance still degrades with many evidence variables

Approximate inference using MCMC

"State" of network = current assignment to all variables

Generate next state by sampling one variable given Markov blanket
 Sample each variable in turn, keeping evidence fixed

```

function MCMC-ASK( $X, e, bn, N$ ) returns an approximation to  $P(X|e)$ 
  local variables:  $N[X]$ , a vector of counts over  $X$ , initially zero
                   $Y$ , the nonevidence variables in  $bn$ 
                   $x$ , the current state of the network, initially copied from  $e$ 

  initialize  $x$  with random values for the variables in  $Y$ 
  for  $j = 1$  to  $N$  do
     $N[x] \leftarrow N[x] + 1$  where  $x$  is the value of  $X$  in  $x$ 
    for each  $Y_i$  in  $Y$  do
      sample the value of  $Y_i$  in  $x$  from  $P(Y_i|MB(Y_i))$  given the values of  $MB(Y_i)$  in  $x$ 
  return NORMALIZE( $N[X]$ )
    
```

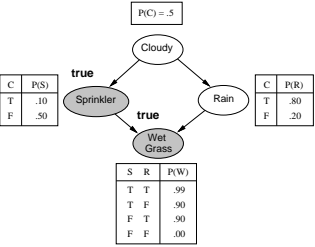
Approaches stationary distribution: long-run fraction of time spent in each state is exactly proportional to its posterior probability

MCMC Example

Estimate $P(Rain|Sprinkler = true, WetGrass = true)$

Sample *Cloudy* then *Rain*, repeat.
 Count number of times *Rain* is true and false in the samples.

Markov blanket of *Cloudy* is *Sprinkler* and *Rain*
 Markov blanket of *Rain* is *Cloudy*, *Sprinkler*, and *WetGrass*



MCMC example contd.

Random initial state: *Cloudy* = true and *Rain* = false

- $P(Cloudy|MB(Cloudy)) = P(Cloudy|Sprinkler, \neg Rain)$
sample \rightarrow false
- $P(Rain|MB(Rain)) = P(Rain|\neg Cloudy, Sprinkler, WetGrass)$
sample \rightarrow true

Visit 100 states
 31 have *Rain* = true, 69 have *Rain* = false

$$\hat{P}(Rain|Sprinkler = true, WetGrass = true) = \text{NORMALIZE}(\langle 31, 69 \rangle) = \langle 0.31, 0.69 \rangle$$

MCMC analysis: Outline

Transition probability $q(\mathbf{y} \rightarrow \mathbf{y}')$

Occupancy probability $\pi_t(\mathbf{y})$ at time t

Equilibrium condition on π_t defines stationary distribution $\pi(\mathbf{y})$
 Note: stationary distribution depends on choice of $q(\mathbf{y} \rightarrow \mathbf{y}')$

Pairwise detailed balance on states guarantees equilibrium

Gibbs sampling transition probability:
 sample each variable given current values of all others
 \Rightarrow detailed balance with the true posterior

For Bayesian networks, Gibbs sampling reduces to sampling conditioned on each variable's Markov blanket

Stationary distribution

$\pi_t(\mathbf{y})$ = probability in state \mathbf{y} at time t
 $\pi_{t+1}(\mathbf{y}')$ = probability in state \mathbf{y}' at time $t + 1$

π_{t+1} in terms of π_t and $q(\mathbf{y} \rightarrow \mathbf{y}')$

$$\pi_{t+1}(\mathbf{y}') = \sum_{\mathbf{y}} \pi_t(\mathbf{y}) q(\mathbf{y} \rightarrow \mathbf{y}')$$

Stationary distribution: $\pi_t = \pi_{t+1} = \pi$

$$\pi(\mathbf{y}') = \sum_{\mathbf{y}} \pi(\mathbf{y}) q(\mathbf{y} \rightarrow \mathbf{y}') \quad \text{for all } \mathbf{y}'$$

If π exists, it is unique (specific to $q(\mathbf{y} \rightarrow \mathbf{y}')$)

In equilibrium, expected "outflow" = expected "inflow"

Detailed balance

"Outflow" = "inflow" for each pair of states:

$$\pi(\mathbf{y}) q(\mathbf{y} \rightarrow \mathbf{y}') = \pi(\mathbf{y}') q(\mathbf{y}' \rightarrow \mathbf{y}) \quad \text{for all } \mathbf{y}, \mathbf{y}'$$

Detailed balance \Rightarrow stationarity:

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

MCMC algorithms typically constructed by designing a transition probability q that is in detailed balance with desired π

Gibbs sampling

Sample each variable in turn, given *all other variables*

Sampling Y_i , let $\bar{\mathbf{Y}}_i$ be all other nonevidence variables

Current values are y_i and $\bar{\mathbf{y}}_i$; \mathbf{e} is fixed

Transition probability is given by

$$q(\mathbf{y} \rightarrow \mathbf{y}') = q(y_i, \bar{\mathbf{y}}_i \rightarrow y'_i, \bar{\mathbf{y}}_i) = P(y'_i | \bar{\mathbf{y}}_i, \mathbf{e})$$

This gives detailed balance with true posterior $P(\mathbf{y} | \mathbf{e})$:

$$\begin{aligned} \pi(\mathbf{y})q(\mathbf{y} \rightarrow \mathbf{y}') &= P(\mathbf{y} | \mathbf{e})P(y'_i | \bar{\mathbf{y}}_i, \mathbf{e}) = P(y_i, \bar{\mathbf{y}}_i | \mathbf{e})P(y'_i | \bar{\mathbf{y}}_i, \mathbf{e}) \\ &= P(y_i | \bar{\mathbf{y}}_i, \mathbf{e})P(\bar{\mathbf{y}}_i | \mathbf{e})P(y'_i | \bar{\mathbf{y}}_i, \mathbf{e}) \quad (\text{chain rule}) \\ &= P(y_i | \bar{\mathbf{y}}_i, \mathbf{e})P(y'_i, \bar{\mathbf{y}}_i | \mathbf{e}) \quad (\text{chain rule backwards}) \\ &= q(\mathbf{y}' \rightarrow \mathbf{y})\pi(\mathbf{y}') = \pi(\mathbf{y}')q(\mathbf{y}' \rightarrow \mathbf{y}) \end{aligned}$$

Markov blanket sampling

A variable is independent of all others given its Markov blanket:

$$P(y'_i | \bar{\mathbf{y}}_i, \mathbf{e}) = P(y'_i | MB(Y_i))$$

Probability given the Markov blanket is calculated as follows:

$$P(y'_i | MB(Y_i)) = P(y'_i | Parents(Y_i)) \prod_{Z_j \in Children(Y_i)} P(z_j | Parents(Z_j))$$

Hence computing the sampling distribution over Y_i for each flip requires just cd multiplications if Y_i has c children and d values; can cache it if c not too large.

Main computational problems:

- 1) Difficult to tell if convergence has been achieved
- 2) Can be wasteful if Markov blanket is large:
 $P(Y_i | MB(Y_i))$ won't change much (law of large numbers)

Performance of approximation algorithms

Absolute approximation: $|P(X | \mathbf{e}) - \hat{P}(X | \mathbf{e})| \leq \epsilon$

Relative approximation: $\frac{|P(X | \mathbf{e}) - \hat{P}(X | \mathbf{e})|}{P(X | \mathbf{e})} \leq \epsilon$

Relative \Rightarrow absolute since $0 \leq P \leq 1$ (may be $O(2^{-n})$)

Randomized algorithms may fail with probability at most δ

Polytime approximation: $\text{poly}(n, \epsilon^{-1}, \log \delta^{-1})$

Theorem (Dagum and Luby, 1993): both absolute and relative approximation for either deterministic or randomized algorithms are NP-hard for any $\epsilon, \delta < 0.5$

(Absolute approximation polytime with no evidence—Chernoff bounds)