

Course Overview

Lecture 10 Inference in Bayesian Networks and Reasoning Over Time

Marco Chiarandini

Department of Mathematics & Computer Science
University of Southern Denmark

Slides by Stuart Russell and Peter Norvig

- ✓ Introduction
 - ✓ Artificial Intelligence
 - ✓ Intelligent Agents
- ✓ Search
 - ✓ Uninformed Search
 - ✓ Heuristic Search
- ✓ Adversarial Search
 - ✓ Minimax search
 - ✓ Alpha-beta pruning
- ✓ Knowledge representation and Reasoning
 - ✓ Propositional logic
 - ✓ First order logic
 - ✓ Inference
- Uncertain knowledge and Reasoning
 - Probability and Bayesian approach
 - Bayesian Networks
 - Hidden Markov Chains
 - Kalman Filters
- Learning
 - Decision Trees
 - Maximum Likelihood
 - EM Algorithm
 - Learning Bayesian Networks
 - Neural Networks
 - Support vector machines

2

Outline

Inference by Random Algs
Exercise
Uncertainty over Time

1. Inference by Randomized Algorithms
2. Exercise
3. Uncertainty over Time

3

Complexity of exact inference

Inference by Random Algs
Exercise
Uncertainty over Time

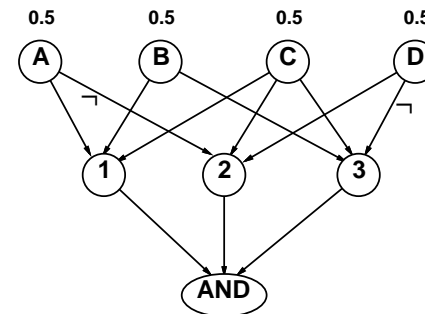
Singly connected networks (or *polytrees*):

- any two nodes are connected by at most one (undirected) path
- time and space cost (with variable elimination) are $O(d^k n)$
- hence time and space cost are linear in n and k bounded by a constant

Multiply connected networks:

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

1. $A \vee B \vee C$
2. $C \vee D \vee \neg A$
3. $B \vee C \vee \neg D$



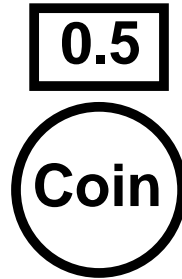
4

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
- Markov chain Monte Carlo (MCMC): sample from a stochastic process whose stationary distribution is the true posterior



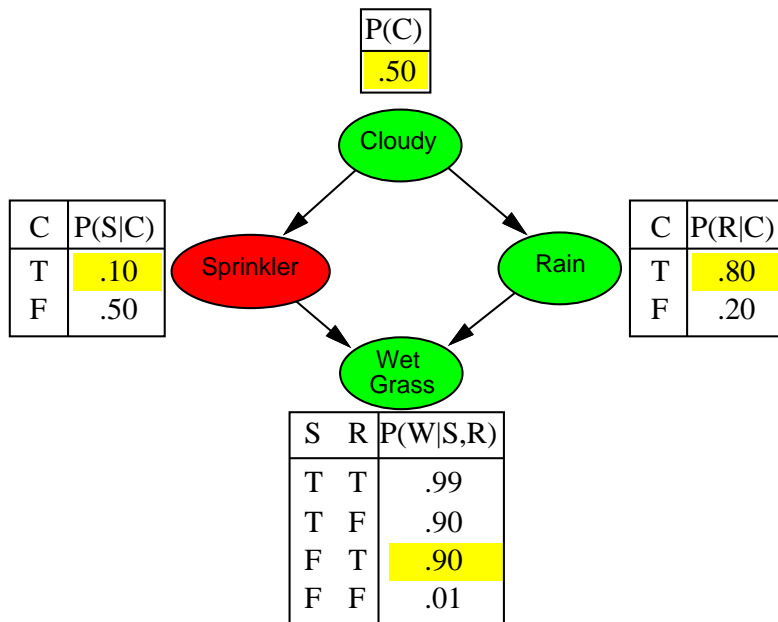
```
function Prior-Sample(bn) returns an event sampled from bn
  inputs: bn, a belief network specifying joint distribution
  P(X1, ..., Xn)

  x ← an event with n elements
  for i = 1 to n do
    xi ← a random sample from P(Xi | parents(Xi))
    given the values of Parents(Xi) in x
  return x
```

5

6

Example



Sampling from an empty network contd.

Probability that `PriorSample` generates a particular event

$$S_{PS}(x_1 \dots x_n) = P(x_1 \dots x_n)$$

i.e., the true prior probability

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

Proof: Let $N_{PS}(x_1 \dots x_n)$ be the number of samples generated for event x_1, \dots, x_n . Then we have

$$\begin{aligned} \lim_{N \rightarrow \infty} \hat{P}(x_1, \dots, x_n) &= \lim_{N \rightarrow \infty} N_{PS}(x_1, \dots, x_n) / N \\ &= S_{PS}(x_1, \dots, x_n) \\ &= \prod_{i=1}^n P(x_i | \text{parents}(X_i)) = P(x_1 \dots x_n) \end{aligned}$$

⇒ That is, estimates derived from `PriorSample` are **consistent**

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

7

8

$\hat{P}(X|e)$ estimated from samples agreeing with 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
     $x \leftarrow$  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( $N[X]$ )
```

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

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

```
function Likelihood-Weighting( $X, e, bn, N$ ) returns an estimate of  $P(X|e)$ 
  local variables:  $W$ , a vector of weighted counts over  $X$ , initially zero

  for  $j = 1$  to  $N$  do
     $x, w \leftarrow$  Weighted-Sample( $bn$ )
     $W[x] \leftarrow W[x] + w$  where  $x$  is the value of  $X$  in  $x$ 
  return Normalize( $W[X]$ )

function Weighted-Sample( $bn, e$ ) 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  $x, w$ 
```

Rejection sampling returns consistent posterior estimates

Proof:

$$\hat{P}(X|e) = \alpha N_{PS}(X, e) \quad (\text{algorithm defn.})$$

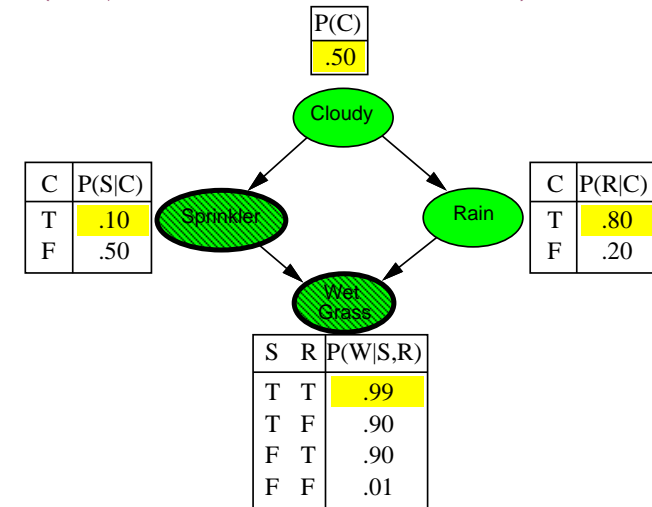
$$= N_{PS}(X, e) / N_{PS}(e) \quad (\text{normalized by } N_{PS}(e))$$

$$\approx P(X, e) / P(e) \quad (\text{property of PriorSample})$$

$$= P(X|e) \quad (\text{defn. of conditional probability})$$

Problem: hopelessly expensive if $P(e)$ is small
 $P(e)$ drops off exponentially with number of evidence variables!

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



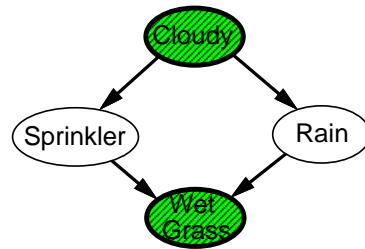
Likelihood weighting analysis

Likelihood weighting returns consistent estimates

Sampling probability for `WeightedSample` is

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

(pays attention to evidence in **ancestors** only)
~somewhere "in between" prior and posterior distribution



Weight for a given sample z, e is

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

but performance still degrades with many evidence variables because a few samples have nearly all the total weight

Weighted sampling probability is

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

13

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 estimate of P(X|e)
  local variables: N[X], a vector of counts over X, initially zero
                  Z, nonevidence variables in bn, hidden + query
                  x, current state of the network, initially copied from e

  initialize x with random values for the variables in Z
  for j = 1 to N do
    N[x] ← N[x] + 1 where x is the value of X in x
    for each Zi in Z do
      sample the value of Zi in x from P(Zi|mb(Zi))
      given the values of MB(Zi) in x
  return Normalize(N[X])
```

Can also choose a variable to sample at random each time

15

Summary

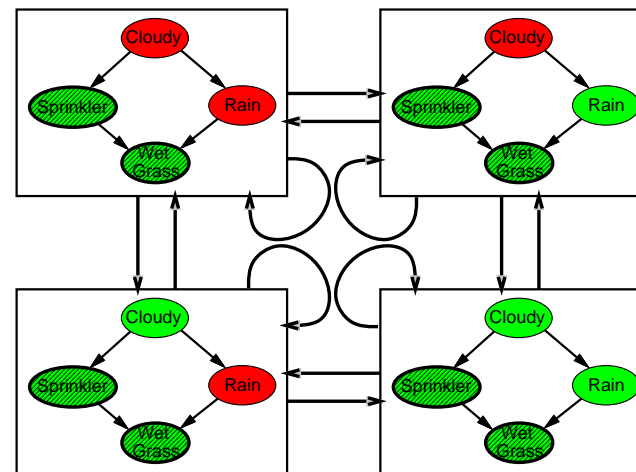
Approximate inference by LW:

- LW does poorly when there is lots of (late-in-the-order) evidence
- LW generally insensitive to topology
- Convergence can be very slow with probabilities close to 1 or 0
- Can handle arbitrary combinations of discrete and continuous variables

14

The Markov chain

With *Sprinkler = true*, *WetGrass = true*, there are four states:



Wander about for a while, average what you see

Probabilistic finite state machine

16

MCMC example contd.

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

Sample *Cloudy* or *Rain* given its Markov blanket, repeat.
Count number of times *Rain* is true and false in the samples.

E.g., visit 100 states

31 have *Rain* = true, 69 have *Rain* = false

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

Theorem

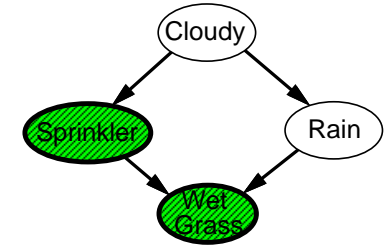
The Markov Chain approaches a *stationary distribution*:
long-run fraction of time spent in each state is exactly
proportional to its posterior probability

17

Markov blanket sampling

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

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



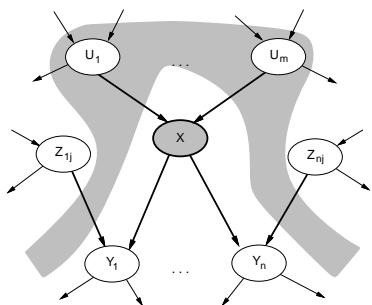
Main computational problems:

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

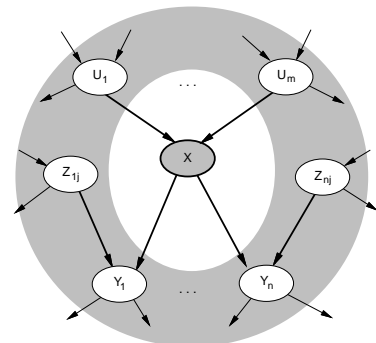
18

Local semantics and Markov Blanket

Local semantics: each node is
conditionally independent
of its nondescendants given its parents



Each node is conditionally
independent of all others given its
Markov blanket: parents + children +
children's parents



19

MCMC analysis: Outline

- Transition probability $q(\mathbf{x} \rightarrow \mathbf{x}')$
- Occupancy probability $\pi_t(\mathbf{x})$ at time t
- Equilibrium condition on π_t defines stationary distribution $\pi(\mathbf{x})$
Note: stationary distribution depends on choice of $q(\mathbf{x} \rightarrow \mathbf{x}')$
- 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

20

Stationary distribution

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

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

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

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

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

- If π exists, it is unique (specific to $q(\mathbf{x} \rightarrow \mathbf{x}')$)
- In equilibrium, expected “outflow” = expected “inflow”

21

Gibbs sampling

- Sample each variable in turn, given **all other variables**
- Sampling X_i , let $\bar{\mathbf{x}}_i$ be all other nonevidence variables
- Current values are x_i and $\bar{\mathbf{x}}_i$; \mathbf{e} is fixed
- Transition probability is given by

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

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

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

23

Detailed balance

- “Outflow” = “inflow” for each pair of states:

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

- Detailed balance \implies stationarity:

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

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

22

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 \implies 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)

24

Exact inference by variable elimination:

- polytime on polytrees, NP-hard on general graphs
- space = time, very sensitive to topology

Approximate inference by LW, MCMC:

- PriorSampling and RejectionSampling unusable as evidence grow
- LW does poorly when there is lots of (late-in-the-order) evidence
- LW, MCMC generally insensitive to topology
- Convergence can be very slow with probabilities close to 1 or 0
- Can handle arbitrary combinations of discrete and continuous variables

1. Inference by Randomized Algorithms

2. Exercise

3. Uncertainty over Time

Wumpus World

Specifying the probability model

1,4	2,4	3,4	4,4
1,3	2,3	3,3	4,3
1,2 B OK	2,2	3,2	4,2
1,1 OK	2,1 B OK	3,1	4,1

$P_{ij} = true$ iff $[i, j]$ contains a pit

$B_{ij} = true$ iff $[i, j]$ is breezy

Include only $B_{1,1}, B_{1,2}, B_{2,1}$ in the probability model

The full joint distribution is $P(P_{1,1}, \dots, P_{4,4}, B_{1,1}, B_{1,2}, B_{2,1})$

Apply product rule: $P(B_{1,1}, B_{1,2}, B_{2,1} | P_{1,1}, \dots, P_{4,4})P(P_{1,1}, \dots, P_{4,4})$

(Do it this way to get $P(Effect|Cause)$.)

First term: 1 if pits are adjacent to breezes, 0 otherwise

Second term: pits are placed randomly, probability 0.2 per square:

$$P(P_{1,1}, \dots, P_{4,4}) = \prod_{i,j=1,1}^{4,4} P(P_{i,j}) = 0.2^n \times 0.8^{16-n}$$

for n pits.

Observations and query

We know the following facts:

$$b = \neg b_{1,1} \wedge b_{1,2} \wedge b_{2,1}$$

$$known = \neg p_{1,1} \wedge \neg p_{1,2} \wedge \neg p_{2,1}$$

Query is $P(P_{1,3} | known, b)$

Define *Unknown* = P_{ij} s other than $P_{1,3}$ and *Known*

For inference by enumeration, we have

$$P(P_{1,3} | known, b) = \alpha \sum_{unknown} P(P_{1,3}, unknown, known, b)$$

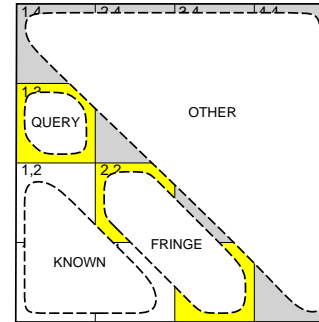
Grows exponentially with number of squares!

Using conditional independence contd.

$$\begin{aligned}
 P(P_{1,3} | known, b) &= \alpha \sum_{unknown} P(P_{1,3}, unknown, known, b) \\
 &= \alpha \sum_{unknown} P(b | P_{1,3}, known, unknown) P(P_{1,3}, known, unknown) \\
 &= \alpha \sum_{fringe} \sum_{other} P(b | known, P_{1,3}, fringe, other) P(P_{1,3}, known, fringe, other) \\
 &= \alpha \sum_{fringe} \sum_{other} P(b | known, P_{1,3}, fringe) P(P_{1,3}, known, fringe, other) \\
 &= \alpha \sum_{fringe} P(b | known, P_{1,3}, fringe) \sum_{other} P(P_{1,3}, known, fringe, other) \\
 &= \alpha \sum_{fringe} P(b | known, P_{1,3}, fringe) \sum_{other} P(P_{1,3}) P(known) P(fringe) P(other) \\
 &= \alpha P(known) P(P_{1,3}) \sum_{fringe} P(b | known, P_{1,3}, fringe) P(fringe) \sum_{other} P(other) \\
 &= \alpha' P(P_{1,3}) \sum_{fringe} P(b | known, P_{1,3}, fringe) P(fringe)
 \end{aligned}$$

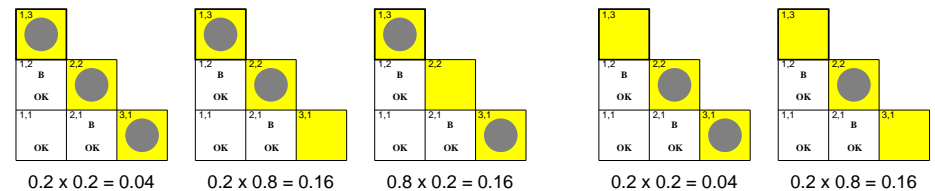
Using conditional independence

Basic insight: observations are conditionally independent of other hidden squares given neighbouring hidden squares



Define $Unknown = Fringe \cup Other$
 $P(b | P_{1,3}, Known, Unknown) = P(b | P_{1,3}, Known, Fringe)$
 Manipulate query into a form where we can use this!

Using conditional independence contd.



$$P(P_{1,3} | known, b) = \alpha' \langle 0.2(0.04 + 0.16 + 0.16), 0.8(0.04 + 0.16) \rangle$$

$$\approx \langle 0.31, 0.69 \rangle$$

$$P(P_{2,2} | known, b) \approx \langle 0.86, 0.14 \rangle$$

1. Inference by Randomized Algorithms
2. Exercise
3. Uncertainty over Time

- ◇ Time and uncertainty
- ◇ Inference: filtering, prediction, smoothing
- ◇ Hidden Markov models
- ◇ Kalman filters (a brief mention)
- ◇ Dynamic Bayesian networks (an even briefer mention)

33

34

Time and uncertainty

- The world changes; we need to track and predict it
- Diabetes management vs vehicle diagnosis
- Basic idea: copy state and evidence variables for each time step
 - \mathbf{X}_t = set of unobservable state variables at time t
e.g., *BloodSugar_t*, *StomachContents_t*, etc.
 - \mathbf{E}_t = set of observable evidence variables at time t
e.g., *MeasuredBloodSugar_t*, *PulseRate_t*, *FoodEaten_t*
- This assumes **discrete time**; step size depends on problem
- Notation: $\mathbf{X}_{a:b} = \mathbf{X}_a, \mathbf{X}_{a+1}, \dots, \mathbf{X}_{b-1}, \mathbf{X}_b$

35

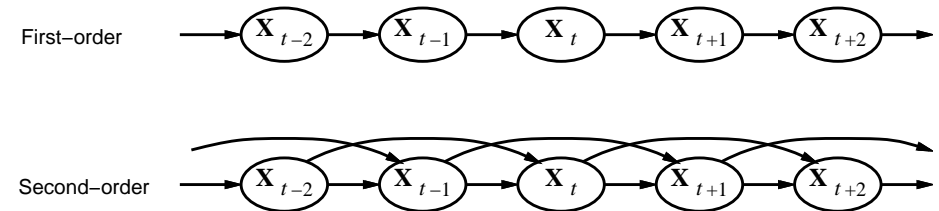
Markov processes (Markov chains)

Construct a Bayes net from these variables:
 - unbounded number of conditional probability table
 - unbounded number of parents

Markov assumption: \mathbf{X}_t depends on **bounded** subset of $\mathbf{X}_{0:t-1}$

First-order Markov process: $P(\mathbf{X}_t | \mathbf{X}_{0:t-1}) = P(\mathbf{X}_t | \mathbf{X}_{t-1})$

Second-order Markov process: $P(\mathbf{X}_t | \mathbf{X}_{0:t-1}) = P(\mathbf{X}_t | \mathbf{X}_{t-2}, \mathbf{X}_{t-1})$



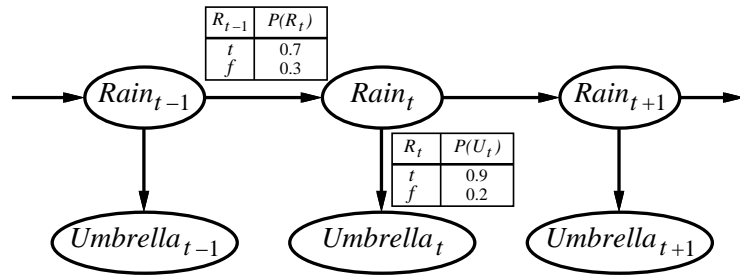
Sensor Markov assumption: $P(\mathbf{E}_t | \mathbf{X}_{0:t}, \mathbf{E}_{0:t-1}) = P(\mathbf{E}_t | \mathbf{X}_t)$

↪ Stationary process:

- transition model $P(\mathbf{X}_t | \mathbf{X}_{t-1})$ and
- sensor model $P(\mathbf{E}_t | \mathbf{X}_t)$ fixed for all t

36

Example



First-order Markov assumption not exactly true in real world!
Possible fixes:

1. **Increase order** of Markov process
2. **Augment state**, e.g., add $Temp_t$, $Pressure_t$

Example: robot motion.

Augment position and velocity with $Battery_t$

Inference tasks

1. **Filtering**: $P(X_t | e_{1:t})$
belief state—input to the decision process of a rational agent
2. **Prediction**: $P(X_{t+k} | e_{1:t})$ for $k > 0$
evaluation of possible action sequences;
like filtering without the evidence
3. **Smoothing**: $P(X_k | e_{1:t})$ for $0 \leq k < t$
better estimate of past states, essential for learning
4. **Most likely explanation**: $\arg \max_{x_{1:t}} P(x_{1:t} | e_{1:t})$
speech recognition, decoding with a noisy channel

37

38

Filtering

Aim: devise a **recursive** state estimation algorithm:

$$P(X_{t+1} | e_{1:t+1}) = f(e_{t+1}, P(X_t | e_{1:t}))$$

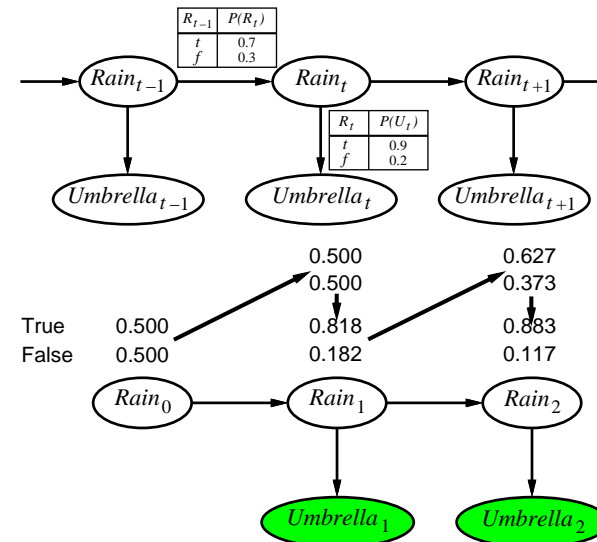
$$\begin{aligned} P(X_{t+1} | e_{1:t+1}) &= P(X_{t+1} | e_{1:t}, e_{t+1}) \\ &= \alpha P(e_{t+1} | X_{t+1}, e_{1:t}) P(X_{t+1} | e_{1:t}) \\ &= \alpha P(e_{t+1} | X_{t+1}) P(X_{t+1} | e_{1:t}) \end{aligned}$$

I.e., **prediction** + **estimation**. Prediction by summing out X_t :

$$\begin{aligned} P(X_{t+1} | e_{1:t+1}) &= \alpha P(e_{t+1} | X_{t+1}) \sum_{x_t} P(X_{t+1} | x_t, e_{1:t}) P(x_t | e_{1:t}) \\ &= \alpha P(e_{t+1} | X_{t+1}) \sum_{x_t} P(X_{t+1} | x_t) P(x_t | e_{1:t}) \end{aligned}$$

$f_{1:t+1} = \text{Forward}(f_{1:t}, e_{t+1})$ where $f_{1:t} = P(X_t | e_{1:t})$
Time and space **constant** (independent of t)

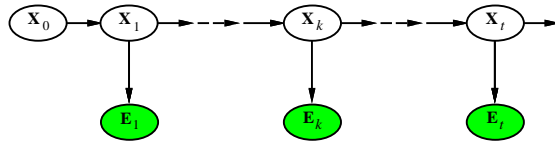
Filtering example



39

40

Smoothing



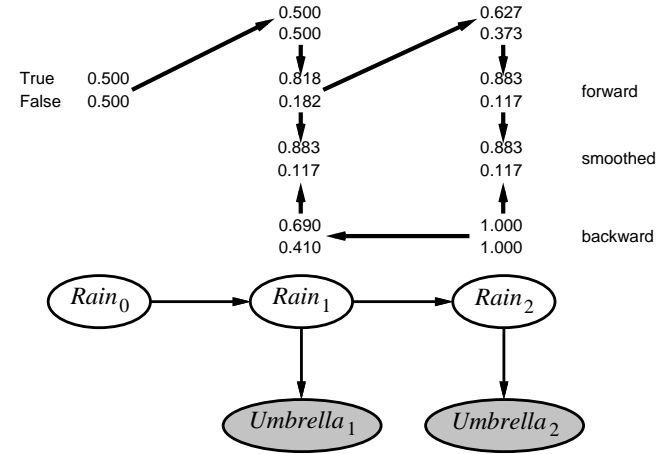
Divide evidence $e_{1:t}$ into $e_{1:k}$, $e_{k+1:t}$:

$$\begin{aligned}
 P(\mathbf{X}_k | \mathbf{e}_{1:t}) &= P(\mathbf{X}_k | \mathbf{e}_{1:k}, \mathbf{e}_{k+1:t}) \\
 &= \alpha P(\mathbf{X}_k | \mathbf{e}_{1:k}) P(\mathbf{e}_{k+1:t} | \mathbf{X}_k, \mathbf{e}_{1:k}) \\
 &= \alpha P(\mathbf{X}_k | \mathbf{e}_{1:k}) P(\mathbf{e}_{k+1:t} | \mathbf{X}_k) \\
 &= \alpha \mathbf{f}_{1:k} \mathbf{b}_{k+1:t}
 \end{aligned}$$

Backward message computed by a backwards recursion:

$$\begin{aligned}
 P(\mathbf{e}_{k+1:t} | \mathbf{X}_k) &= \sum_{x_{k+1}} P(\mathbf{e}_{k+1:t} | \mathbf{X}_k, x_{k+1}) P(x_{k+1} | \mathbf{X}_k) \\
 &= \sum_{x_{k+1}} P(\mathbf{e}_{k+1:t} | x_{k+1}) P(x_{k+1} | \mathbf{X}_k) \\
 &= \sum_{x_{k+1}} P(\mathbf{e}_{k+1:t} | x_{k+1}) P(\mathbf{e}_{k+2:t} | x_{k+1}) P(x_{k+1} | \mathbf{X}_k)
 \end{aligned}$$

Smoothing example



Forward-backward algorithm: cache forward messages along the way
Time linear in t (polytree inference), space $O(t|f)$