Inference in Bayes Nets AI Class 10 (Ch. 14.1—14.4.2; skim 14.3)





How is the Bayesian network created?

- 1. Get an expert to design it
 - Expert must determine the structure of the Bayesian network
 - This is best done by modeling direct causes of a variable as its parents
 - Expert must determine the values of the CPT entries
 - These values could come from the expert's informed opinion
 - Or an external source like census information
 - Or they are estimated from data
 - Or a combination of the above
- 2. Learn it from data
 - This is a much better option but it usually requires a large amount of data
 - This is where Bayesian statistics comes in!



Probability, redux

- Worlds, random variables, events, sample space
- Joint probabilities of multiple connected variables
- Conditional probabilities of a variable, given another variable(s)
- Marginalizing out unwanted variables
- Inference from the joint probability

The big idea: figuring out the probability of variable(s) taking certain value(s)

Review: Independence

What does it mean for A and B to be independent?

- P(A) **L** P(B)
- A and B do not affect each other's probability
- $P(A \land B) = P(A) P(B)$

6

Review: Conditioning

What does it mean for A and B to be **conditionally independent given C?**

- A and B don't affect each other if C is known
- $P(A \land B \mid C) = P(A \mid C) P(B \mid C)$



Review: Bayes' Rule What is Bayes' Rule? $P(H_i | E_j) = \frac{P(E_j | H_i)P(H_i)}{P(E_j)}$ What's it useful for? • Diagnosis • Effect is perceived, want to know (probability of) cause $P(hidden | observed) = \frac{P(observed | hidden)P(hidden)}{P(observed)}$

R&N, 495–496

Review: Joint Probability

- What is the joint probability of A and B?
 - *P*(A,B)
- The probability of **any pair** of legal assignments.
 - Generalizing to > 2, of course
- Booleans: expressed as a matrix/table

	alarm	¬ alarm	
burglary	0.09	0.01	
¬ burglary	0.1	0.8	

Α	В	
Т	Т	0.09
Т	F	0.1
F	Т	0.01
F	F	0.8

Continuous domains: probability functions

10

Next Up **Bayesian networks** Network structure and independence • Inference in Bayesian networks • Exact inference В D Approximate inference C E Ε -G G DAG: In a Bayes net, the links may form loops, but they may not form cycles.

Review: Bayes' Nets: Big Picture

- Problems with full joint distribution tables as our probabilistic models:
 - Joint gets way too big to represent explicitly
 - Unless there are only a few variables
 - Hard to learn (estimate) anything empirically about more than a few variables at a time

	Α		¬A	
	Е	¬E	Е	¬Ε
В	0.01	0.08	0.001	0.009
¬₿	0.01	0.09	0.01	0.79

Slides derived from Matt E. Taylor, U Alberta

12

Review: Bayes' Nets Bayesian Network: **BN = (DAG, CPD)** DAG: directed acyclic graph (BN's structure) **CPD**: conditional probability distribution (BN's parameters) P(A) = 0.001P(B|A) = 0.3P(C|A) = 0.2 $P(B|\neg A) = 0.001$ $P(C|\neg A) = 0.005$ (В) C $P(\neg B|A) = 0.7$ $P(\neg B | \neg A) = 0.999$ (E) P(E|C) = 0.4D $P(E|\neg C) = 0.002$ P(D|B,C) = 0.1 $P(D|B,\neg C) = 0.01$ $P(D|\neg B,C) = 0.01$ $P(D|\neg B,\neg C) = 0.00001$



The Chain Rule

- $P(\alpha_1 \land \alpha_2 \land ... \land \alpha_n) = P(\alpha_1) \times$ $P(\alpha_2 \mid \alpha_1) \times$ $P(\alpha_3 \mid \alpha_1 \land \alpha_2) \times ... \times$ $P(\alpha_n \mid \alpha_1 \land \cdots \land \alpha_{n-1})$
 - $= \prod_{i=1..n} P(\alpha_i \mid \alpha_1 \land \dots \land \alpha_{i-1})$
 - $= P(x_1,...,x_n) = \prod_{i=1}^n P(x_i \mid \pi_i)$

artint.info/html/ArtInt_143.html







5 A node is conditionally independent of its non-descendants given its parents. 6 A node is conditionally independent of all other nodes in the network given its parents, children, and children's parents (also known as its Markov blanket)

The Joint Probability Distribution

• Due to the Markov condition, we can compute the joint probability distribution over all the variables X₁, ..., X_n in the Bayesian net using the formula:

$$P(X_1 = x_1, ..., X_n = x_n) = \prod_{i=1}^n P(X_i = x_i \mid Parents(X_i))$$

Where $Parents(X_i)$ means the values of the Parents of the node X_i with respect to the graph

21

Independence and Causal Chains

- Important question about a BN:
 - Are two nodes independent given certain evidence?
 - If yes, we can it prove using algebra (tedious)
 - If no, can prove it with a counter-example
- Question: are X and Z necessarily independent?
 - No.
 - Ex: Clouds (X) cause rain (Y), which causes traffic (Z)
 - X can influence Z, Z can influence X (via Y)
- This configuration is a "causal chain"













Representational Extensions

- Conditional probability tables (CPTs) for large networks can require a large number of parameters
 - O(2^k) where k is the branching factor of the network
- There are ways of compactly representing CPTs
 - Deterministic relationships
 - Noisy-OR
 - Noisy-MAX
- What about continuous variables?
 - Discretization
 - Use density functions (usually mixtures of Gaussians) to build hybrid Bayesian networks (with discrete and continuous variables)





What Are We Trying To Do?

- Now we know the semantics of Bayes' nets
- But how do we use it?
- Say we have some evidence (that is, some variables are instantiated)
- We usually want to know the probability of some **other** variables
- Why?
 - Reason about hidden (non-observed) information
 - What caused something?
 - What is the **probability** of something?



Inference

- Instead of computing the joint, suppose we just want the probability for one variable (or a subset)
- Using a Bayesian network to compute probabilities is called inference
- In general, inference involves queries of the form:



Inference Techniques

- **Exact inference:** Analytically compute the conditional probability distribution over the variables we care about
- Approximate inference: Sometimes exact inference is too hard
 - Come up with approximate solutions based on statistical sampling

Exact inference

- Enumeration
- Belief propagation in polytrees
- Variable elimination
- Clustering / join tree algorithms

Approximate inference

- Stochastic simulation / sampling methods
- Markov chain Monte Carlo methods
- Genetic algorithms
- Neural networks
- Simulated annealing
- Mean field theory

Query Types

Given a Bayesian network, what questions might we want to ask?

- Conditional probability query: P(x | e)
 - Given instantiations for some of the variables (e = values of all instantiated variables; it doesn't have to be just one), what is the probability that node X has a particular value x?

Query Types

Given a Bayesian network, what questions might we want to ask?

- Conditional probability query: P(x | e)
 - Given instantiations for some of the variables (e = values of all instantiated variables; it doesn't have to be just one), what is the probability that node X has a particular value x?
- Maximum a posteriori probability: What value of x maximizes P(x|e)?
 - What is the most likely explanation for some evidence?
 - That is, what is the value of node(s) X that maximizes the probability that you would have seen the evidence you did?
 - This is called a maximum a posteriori probability or MAP query



38

Inference Tasks

- Simple queries: Compute posterior marginal P(X_i | E=value)
 - E.g., P(NoGas | Gauge=*empty*, Lights=*on*, Starts=*false*)
- Conjunctive queries:
 - $P(X_i, X_j | E=value) = P(X_i | E=value) P(X_j | X_i, E=value)$
- Optimal decisions:
 - *Decision networks* include utility information
 - Probabilistic inference gives P(outcome | action, evidence)
- Value of information: Which evidence should we seek next?
- Sensitivity analysis: Which probability values are most critical?
- Explanation: Why do I need a new starter motor?

Using the joint distribution

• To answer any query involving a conjunction of variables, sum over the variables not involved in the query

40

Using the joint distribution • To answer any query involving a conjunction of variables, sum over the variables not involved in the query $Pr(d) = \sum_{ABC} Pr(a,b,c,d)$ $= \sum_{ABC} \sum_{ABC} \sum_{ABC} \sum_{C} Pr(A = a \land B = b \land C = c)$

Using the joint distribution

• To answer any query involving a conjunction of variables, sum over the variables not involved in the query

$$\Pr(d) = \sum_{ABC} \Pr(a, b, c, d)$$
$$= \sum_{a \in \operatorname{dom}(A)} \sum_{b \in \operatorname{dom}(B)} \sum_{c \in \operatorname{dom}(C)} \Pr(A = a \land B = b \land C = c)$$

42

Using the joint distribution • To answer any query involving a conjunction of variables, sum over the variables not involved in the query $Pr(d) = \sum_{ABC} Pr(a,b,c,d)$ $= \sum_{a \in dom(A)b \in dom(B)c \in dom(C)} Pr(A = a \land B = b \land C = c)$ Summing over A and C, because b and d are instantiated in the query $Pr(d \mid b) = \frac{Pr(b,d)}{Pr(b)} = \sum_{ACD} Pr(a,b,c,d)$

Inference by Enumeration

- Add all of the terms (atomic event probabilities) from the full joint distribution
- If E are the evidence (observed) variables and Y are the other (unobserved) variables, then:
 - $P(X | E) = \alpha P(X, E) = \alpha \sum P(X, E, Y)$
- Each P(X, E, Y) term can be computed using the chain rule
- Computationally expensive!





Reminder: P(E) is known (observed), so 1/P(E) is a constant that makes everything sum to 1: the *normalizing constant*















Simple Case $A \longrightarrow B \longrightarrow C \longrightarrow D$ $Pr(d) = \sum_{ABC} Pr(a,b,c,d)$ $= \sum_{ABC} Pr(d \mid c) Pr(c \mid b) Pr(b \mid a) Pr(a)$











Simple Case $A \xrightarrow{B} \xrightarrow{C} \xrightarrow{D} D$ $Pr(d) = \sum_{C} Pr(d | c) \sum_{B} Pr(c | b) \sum_{A} Pr(b | a) Pr(a)$ $f_{1}(b)$













Variable Elimination Algorithm

• Given a Bayesian network, and an *elimination order* for the non-query variables























































Properties of Variable Elimination

- Time is exponential in size of largest factor
- Bad elimination order can generate huge factors
- NP Hard to find the best elimination order
- Even the best elimination order may generate large factors
- There are reasonable heuristics for picking an elimination order (such as choosing the variable that results in the smallest next factor)
- Inference in polytrees (nets with no cycles) is linear in size of the network (the largest CPT)























Learning Bayesian Networks from Data

- We won't have enough time to describe how we actually learn Bayesian networks from data ☺
- Some references:
 - Gregory F. Cooper and Edward Herskovits. A Bayesian Method for the Induction of Probabilistic Networks from Data. Machine Learning, 9:309-347, 1992.
 - David Heckerman. A Tutorial on Learning Bayesian Networks. Technical Report MSR-TR-95-06, Microsoft Research. 1995. (Available online)

Other Inference Methods

- Convert network to a polytree
 - In a polytree no two nodes have more than one path between them
 - For such a graph there is a linear time algorithm
 - However, converting into a polytree requires a large increase in the size of the graph (number of nodes)
- Why is inference in polytrees easy?
- Given a variable X we can always divide the other variables into two sets:
 - E+: Variables 'above' X
 - E-: Variables 'below' X
- These sets are mutually exclusive
- Using these sets we can efficiently compute conditional and joint distributions

114

Summary

- Bayes nets
 - Structure
 - Parameters
 - Conditional independence
 - Chaining
- BN inference
 - Enumeration
 - Variable elimination
 - Sampling methods..?













