

Statistical Reasoning

Probability and Bayes' theorem and causal networks, reasoning belief network

Introduction



Suppose you are trying to determine if a patient has inhalational anthrax. You observe the following symptoms:

- The patient has a cough
- The patient has a fever
- The patient has difficulty breathing

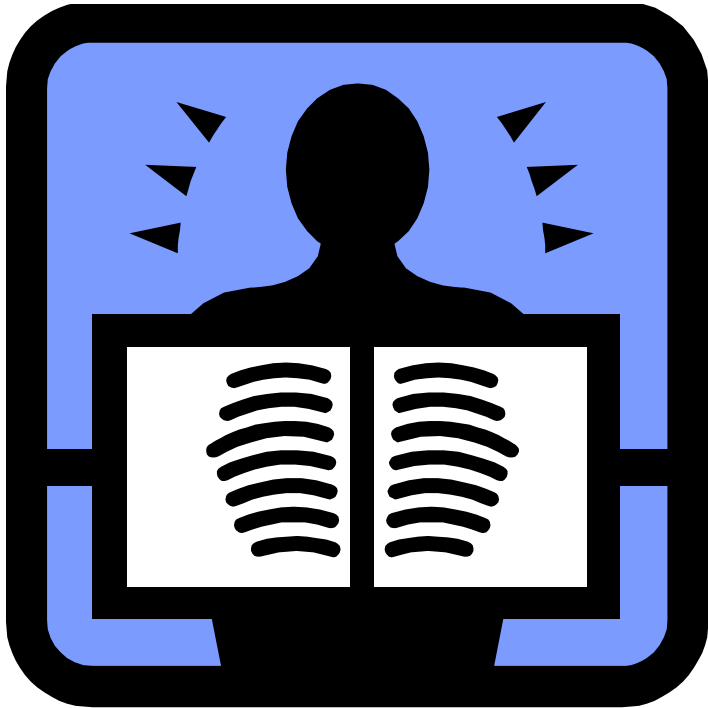
Introduction



You would like to determine how likely the patient is infected with inhalational anthrax given that the patient has a cough, a fever, and difficulty breathing

We are not 100% certain that the patient has anthrax because of these symptoms. We are dealing with uncertainty!

Introduction



Now suppose you order an x-ray and observe that the patient has a wide mediastinum.

Your belief that that the patient is infected with inhalational anthrax is now much higher.

Introduction

- In the previous slides, what you observed affected your belief that the patient is infected with anthrax
- This is called reasoning with uncertainty
- Wouldn't it be nice if we had some methodology for reasoning with uncertainty? Why in fact, we do...

- How does these uncertainty come??

Sources of Uncertainty

- Uncertain **inputs** -- missing and/or noisy data
- Uncertain **knowledge**
 - Multiple causes lead to multiple effects
 - Incomplete enumeration of conditions or effects
 - Incomplete knowledge of causality in the domain
 - Probabilistic/stochastic effects
- Uncertain **outputs**
 - Abduction and induction are inherently uncertain
 - Default reasoning, even deductive, is uncertain
 - Incomplete deductive inference may be uncertain
- ▶ Probabilistic reasoning only gives probabilistic results
(summarizes uncertainty from various sources)

Decision making with uncertainty

Rational behavior:

- For each possible action, identify the possible outcomes
- Compute the **probability** of each outcome
- Compute the **utility** of each outcome
- Compute the probability-weighted **(expected) utility** over possible outcomes for each action
- Select action with the highest expected utility (principle of **Maximum Expected Utility**)

At a glance

- if we roll two dice, each showing one of six possible numbers, the number of total unique rolls is $6*6 = 36$. We distinguish the dice in some way (a first and second or left and right die). Here is a listing of the joint possibilities for the dice:

(1,1) (1,2) (1,3) (1,4) (1,5) (1,6)
(2,1) (2,2) (2,3) (2,4) (2,5) (2,6)
(3,1) (3,2) (3,3) (3,4) (3,5) (3,6)
(4,1) (4,2) (4,3) (4,4) (4,5) (4,6)
(5,1) (5,2) (5,3) (5,4) (5,5) (5,6)
(6,1) (6,2) (6,3) (6,4) (6,5) (6,6)

- The number of rolls which add up to 4 is 3 ((1,3), (2,2), (3,1)), so the probability of rolling a total of 4 is $3/36 = 1/12$.
- This does not mean 8.3% true, but 8.3% chance of it being true.

Probabilities anyway?

Kolmogorov showed that three simple axioms lead to the rules of probability theory

1. All probabilities are between 0 and 1:

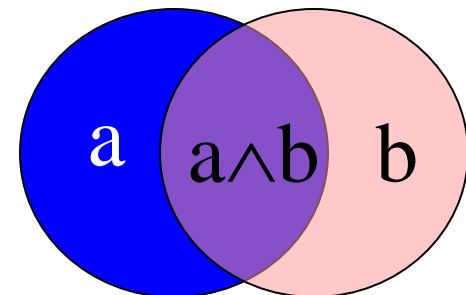
$$0 \leq P(a) \leq 1$$

2. Valid propositions (tautologies) have probability 1, and unsatisfiable propositions have probability 0:

$$P(\text{true}) = 1 ; P(\text{false}) = 0$$

3. The probability of a disjunction is given by:

$$P(a \vee b) = P(a) + P(b) - P(a \wedge b)$$



Probability theory

- **Random variables**

 - Domain

- **Atomic event**: complete specification of state

- **Prior probability**: degree of belief without any other evidence

- **Joint probability**: matrix of combined probabilities of a set of variables

Probability theory

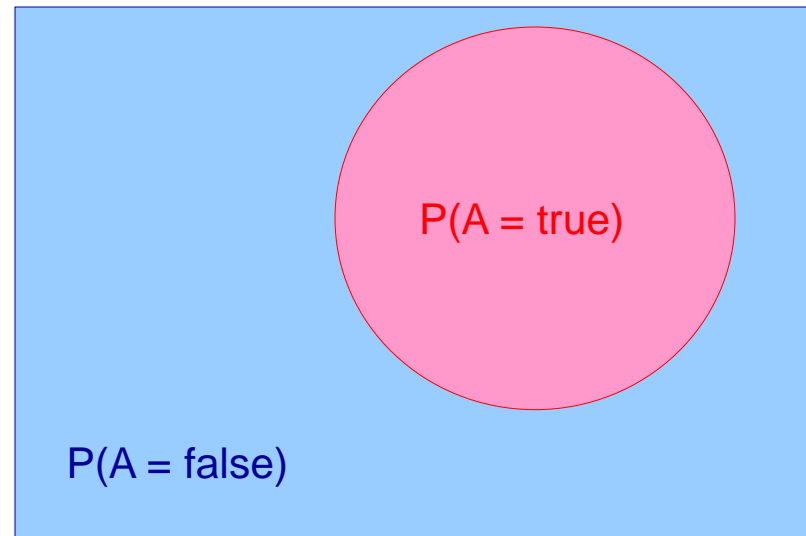
- **Conditional probability**: prob. of effect given causes
- **Computing conditional probs**:
 - $P(a \mid b) = P(a \wedge b) / P(b)$
 - $P(b)$: **normalizing** constant
- **Product rule**:
 - $P(a \wedge b) = P(a \mid b) * P(b)$

Probabilities Theory

We will write $P(A = \text{true})$ to mean the probability that $A = \text{true}$.

What is probability? It is the relative frequency with which an outcome would be obtained if the process were repeated a large number of times under similar conditions*

The sum of the red
and blue areas is 1

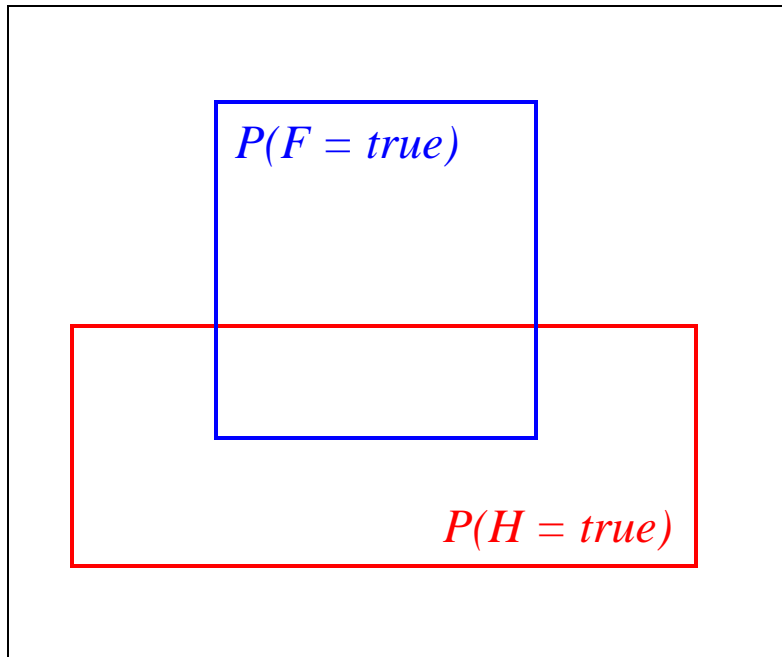


*Ahem...there's also the Bayesian definition which says probability is your degree of belief in an outcome



Conditional Probability

- $P(A = \text{true} \mid B = \text{true})$ = Out of all the outcomes in which B is true, how many also have A equal to true
- Read this as: “Probability of A conditioned on B ” or “Probability of A given B ”



H = “Have a headache”

F = “Coming down with Flu”

$$P(H = \text{true}) = 1/10$$

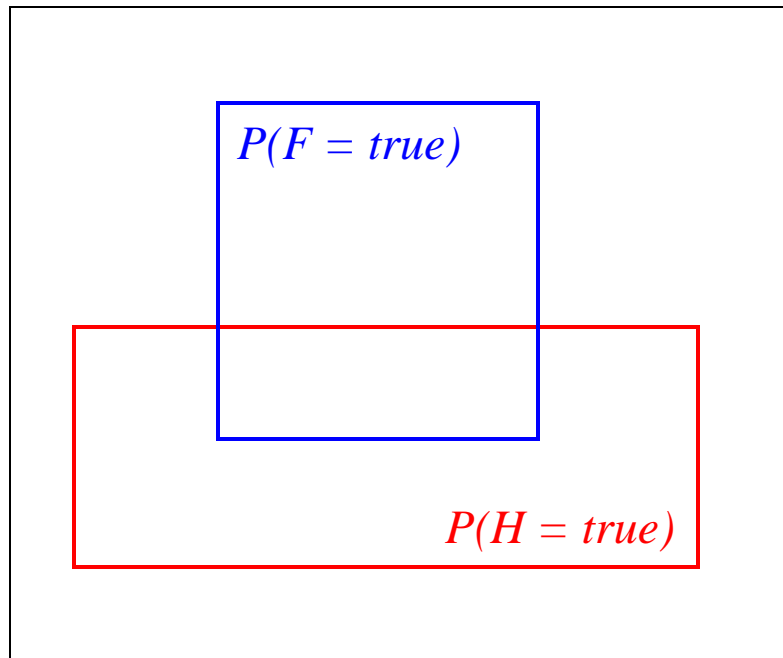
$$P(F = \text{true}) = 1/40$$

$$P(H = \text{true} \mid F = \text{true}) = 1/2$$

“Headaches are rare and flu is rarer, but if you’re coming down with flu there’s a 50-50 chance you’ll have a headache.”

The Joint Probability Distribution

- We will write $P(A = \text{true}, B = \text{true})$ to mean “the probability of $A = \text{true}$ **and** $B = \text{true}$ ”
- Notice that:



$$\begin{aligned} &P(H = \text{true} | F = \text{true}) \\ &= \frac{\text{Area of "H and F" region}}{\text{Area of "F" region}} \\ &= \frac{P(H = \text{true}, F = \text{true})}{P(F = \text{true})} \end{aligned}$$

In general, $P(X|Y) = P(X, Y) / P(Y)$

The Joint Probability Distribution

- Joint probabilities can be between any number of variables
eg. $P(A = \text{true}, B = \text{true}, C = \text{true})$
- For each combination of variables, we need to say how probable that combination is
- The probabilities of these combinations need to sum to 1

A	B	C	P(A,B,C)
false	false	false	0.1
false	false	true	0.2
false	true	false	0.05
false	true	true	0.05
true	false	false	0.3
true	false	true	0.1
true	true	false	0.05
true	true	true	0.15

Sums to 1

The Joint Probability Distribution

- Once you have the joint probability distribution, you can calculate any probability involving A , B , and C
- Note: May need to use marginalization and Bayes rule, (both of which are not discussed in these slides)

A	B	C	P(A,B,C)
false	false	false	0.1
false	false	true	0.2
false	true	false	0.05
false	true	true	0.05
true	false	false	0.3
true	false	true	0.1
true	true	false	0.05
true	true	true	0.15

Examples of things you can compute:

- $P(A=true) = \text{sum of } P(A,B,C) \text{ in rows with } A=true$
- $P(A=true, B = true \mid C=true) =$
 $P(A = true, B = true, C = true) / P(C = true)$

The Problem with the Joint Distribution

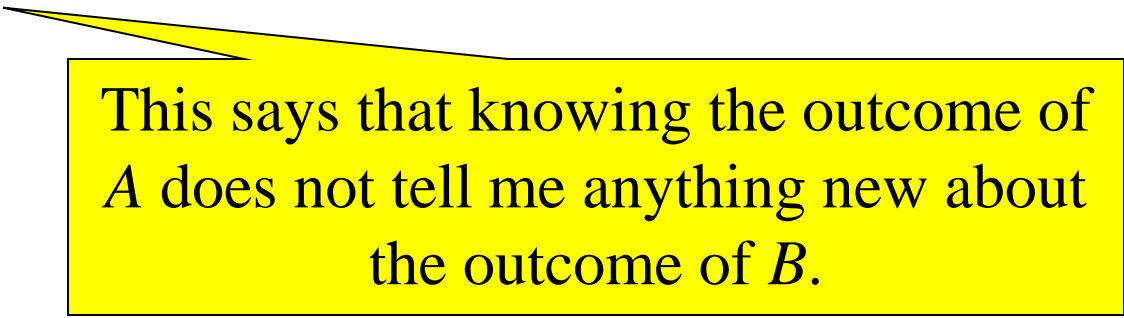
- Lots of entries in the table to fill up!
- For k Boolean random variables, you need a table of size 2^k
- How do we use fewer numbers? Need the concept of independence

A	B	C	P(A,B,C)
false	false	false	0.1
false	false	true	0.2
false	true	false	0.05
false	true	true	0.05
true	false	false	0.3
true	false	true	0.1
true	true	false	0.05
true	true	true	0.15

Independence

Variables A and B are independent if any of the following hold:

- $P(A, B) = P(A) P(B)$
- $P(A \mid B) = P(A)$
- $P(B \mid A) = P(B)$



This says that knowing the outcome of A does not tell me anything new about the outcome of B .

Independence

How is independence useful?

- Suppose you have n coin flips and you want to calculate the joint distribution $\mathbf{P}(C_1, \dots, C_n)$
- If the coin flips are not independent, you need 2^n values in the table
- If the coin flips are independent, then

$$P(C_1, \dots, C_n) = \prod_{i=1}^n P(C_i)$$

Each $P(C_i)$ table has 2 entries and there are n of them for a total of $2n$ values

Conditional Independence

Variables A and B are conditionally independent given C if any of the following hold:

- $P(A, B \mid C) = P(A \mid C) P(B \mid C)$
- $P(A \mid B, C) = P(A \mid C)$
- $P(B \mid A, C) = P(B \mid C)$

Knowing C tells me everything about B . I don't gain anything by knowing A (either because A doesn't influence B or because knowing C provides all the information knowing A would give)



Independence

- When sets of variables don't affect each others' probabilities, we call them **independent**, and can easily compute their joint and conditional probability:

$$\text{Independent}(A, B) \rightarrow P(A \wedge B) = P(A) * P(B), P(A | B) = P(A)$$

- {moonPhase, lightLevel} *might* be independent of {burglary, alarm, earthquake}
 - Maybe not: crooks may be more likely to burglarize houses during a new moon (and hence little light)
 - But if we know the light level, the moon phase doesn't affect whether we are burglarized
 - If burglarized, light level doesn't affect if alarm goes off
- Need a more complex notion of independence and methods for reasoning about the relationships

Axioms of Probability

- Bayes' Rule
 - Given a hypothesis (H) and evidence (E), and given that $P(E) > 0$, what is $P(H|E)$?
- Many times rules and information are uncertain, yet we still want to say something about the consequent; namely, the degree to which it can be believed. A British cleric and mathematician, Thomas Bayes, suggested an approach.
- Recall the two forms of the product rule:
 - $P(ab) = P(a) * P(b|a)$
 - $P(ab) = P(b) * P(a|b)$
- If we equate the two right-hand sides and divide by $P(a)$, we get
$$P(b|a) = \frac{P(a|b)P(b)}{P(a)}$$

Example

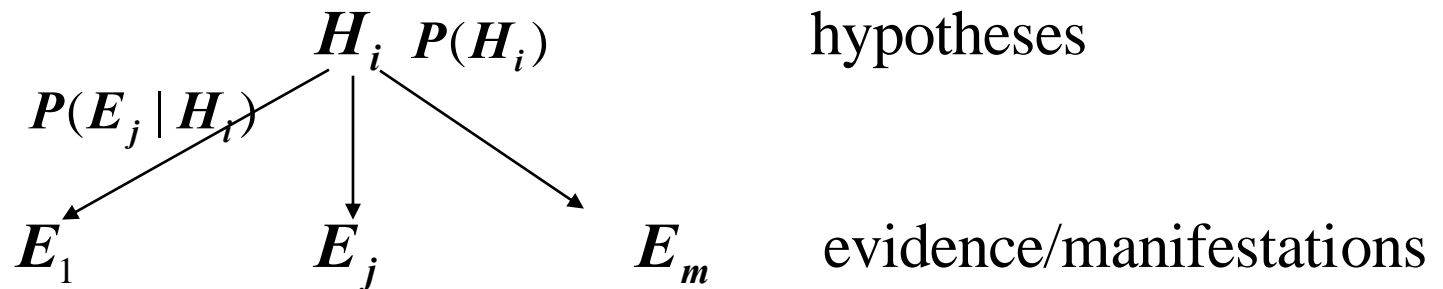
- Bayes' rule is useful when we have three of the four parts of the equation.
- For example, a doctor knows that meningitis causes a stiff neck in 50% of such cases. The prior probability of having meningitis is 1/50,000 and the prior probability of any patient having a stiff neck is 1/20.
- What is the probability that a patient has meningitis if they have a stiff neck?
- H = "Patient has meningitis"
- E = "Patient has stiff neck"

$$P(H|E) = \frac{P(E|H) * P(H)}{P(E)}$$

$$P(H|E) = (0.5 * .00002) / .05 = .0002$$

Bayesian inference

- In the setting of diagnostic/evidential reasoning



- Know prior probability of hypothesis
- conditional probability

$$P(H_i)$$

$$P(E_j | H_i)$$

$$P(H_i | E_j)$$

- Want to compute the *posterior probability*

- Bayes' s theorem (formula 1):

$$P(H_i | E_j) = P(H_i) * P(E_j | H_i) / P(E_j)$$

Simple Bayesian diagnostic reasoning

- Also known as: [Naive Bayes classifier](#)
- Knowledge base:
 - Evidence / manifestations: E_1, \dots, E_m
 - Hypotheses / disorders: H_1, \dots, H_n
 - Note: E_j and H_i are **binary**; hypotheses are **mutually exclusive** (non-overlapping) and **exhaustive** (cover all possible cases)
 - Conditional probabilities: $P(E_j \mid H_i), i = 1, \dots, n; j = 1, \dots, m$
- Cases (evidence for a particular instance): E_1, \dots, E_l
- Goal: Find the hypothesis H_i with the highest posterior
 - $\text{Max}_i P(H_i \mid E_1, \dots, E_l)$

Simple Bayesian diagnostic reasoning

- Bayes' rule says that

$$P(H_i \mid E_1 \dots E_m) = P(E_1 \dots E_m \mid H_i) P(H_i) / P(E_1 \dots E_m)$$

- Assume each evidence E_i is conditionally independent of the others, *given* a hypothesis H_i , then:

$$P(E_1 \dots E_m \mid H_i) = \prod_{j=1}^m P(E_j \mid H_i)$$

- If we only care about relative probabilities for the H_i , then we have:

$$P(H_i \mid E_1 \dots E_m) = \alpha P(H_i) \prod_{j=1}^m P(E_j \mid H_i)$$

Limitations

- Cannot easily handle multi-fault situations, nor cases where intermediate (hidden) causes exist:
 - Disease D causes syndrome S, which causes correlated manifestations M_1 and M_2
- Consider a composite hypothesis $H_1 \wedge H_2$, where H_1 and H_2 are independent. What's the relative posterior?

$$\begin{aligned} P(H_1 \wedge H_2 \mid E_1, \dots, E_l) &= \alpha P(E_1, \dots, E_l \mid H_1 \wedge H_2) P(H_1 \wedge H_2) \\ &= \alpha P(E_1, \dots, E_l \mid H_1 \wedge H_2) P(H_1) P(H_2) \\ &= \alpha \prod_{j=1}^l P(E_j \mid H_1 \wedge H_2) P(H_1) P(H_2) \end{aligned}$$

- How do we compute $P(E_j \mid H_1 \wedge H_2)$?

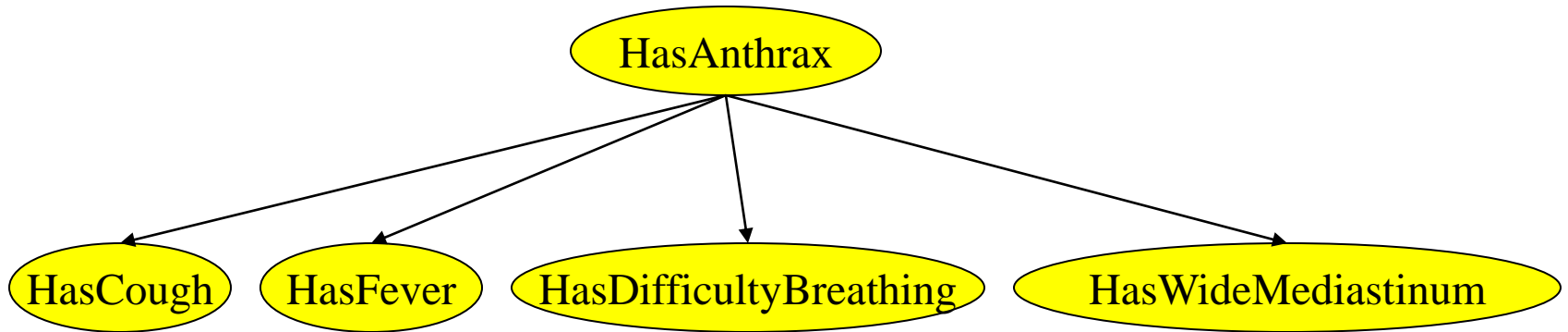
Limitations

- Assume H_1 and H_2 are independent, given E_1, \dots, E_l ?
 - $P(H_1 \wedge H_2 \mid E_1, \dots, E_l) = P(H_1 \mid E_1, \dots, E_l) P(H_2 \mid E_1, \dots, E_l)$
- This is a very unreasonable assumption
 - Earthquake and Burglar are independent, but *not* given Alarm:
 - $P(\text{burglar} \mid \text{alarm}, \text{earthquake}) \ll P(\text{burglar} \mid \text{alarm})$
- Another limitation is that simple application of Bayes' s rule doesn't allow us to handle causal chaining:
 - A: this year's weather; B: cotton production; C: next year's cotton price
 - A influences C indirectly: $A \rightarrow B \rightarrow C$
 - $P(C \mid B, A) = P(C \mid B)$
- Need a richer representation to model interacting hypotheses, conditional independence, and causal chaining
- Next: conditional independence and Bayesian networks!

Summary

- Probability is a rigorous formalism for uncertain knowledge
- **Joint probability distribution** specifies probability of every atomic event
- Can answer queries by summing over atomic events
- But we must find a way to reduce the joint size for non-trivial domains
- **Bayes' rule** lets unknown probabilities be computed from known conditional probabilities, usually in the causal direction
- **Independence** and **conditional independence** provide the tools

Bayesian Networks

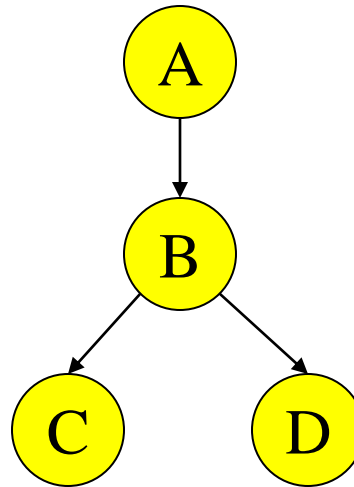


- In the opinion of many AI researchers, Bayesian networks are the most significant contribution in AI in the last 10 years
- They are used in many applications eg. spam filtering, speech recognition, robotics, diagnostic systems and even syndromic surveillance

A Bayesian Network

A Bayesian network is made up of:

1. A Directed Acyclic Graph



2. A set of tables for each node in the graph

A	P(A)
false	0.6
true	0.4

A	B	P(B A)
false	false	0.01
false	true	0.99
true	false	0.7
true	true	0.3

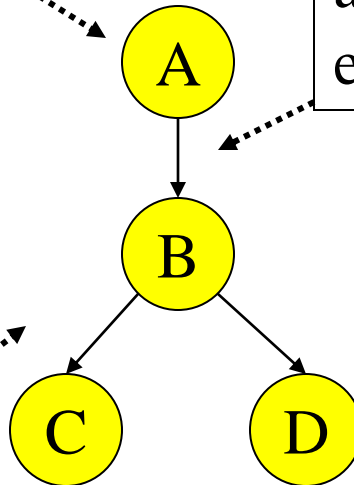
B	D	P(D B)
false	false	0.02
false	true	0.98
true	false	0.05
true	true	0.95

B	C	P(C B)
false	false	0.4
false	true	0.6
true	false	0.9
true	true	0.1

A Directed Acyclic Graph

Each node in the graph is a random variable

A node X is a parent of another node Y if there is an arrow from node X to node Y
eg. A is a parent of B



Informally, an arrow from node X to node Y means X has a direct influence on Y

A Set of Tables for Each Node

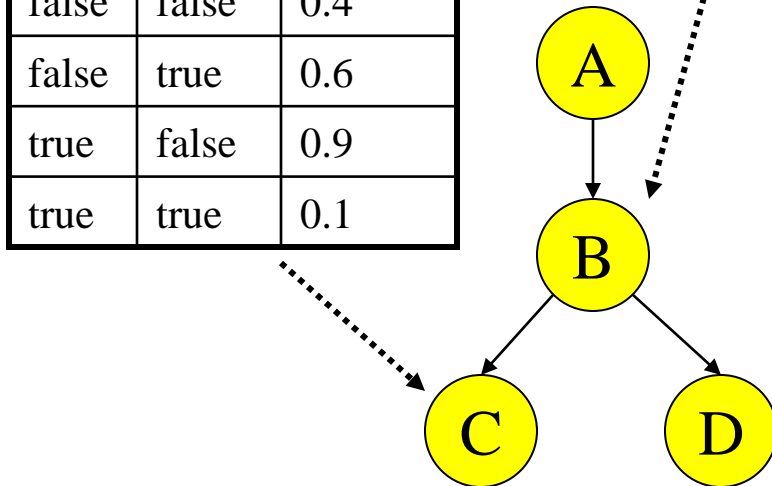
A	P(A)
false	0.6
true	0.4

A	B	P(B A)
false	false	0.01
false	true	0.99
true	false	0.7
true	true	0.3

B	C	P(C B)
false	false	0.4
false	true	0.6
true	false	0.9
true	true	0.1

Each node X_i has a conditional probability distribution $P(X_i \mid \text{Parents}(X_i))$ that quantifies the effect of the parents on the node

The parameters are the probabilities in these conditional probability tables (CPTs)



B	D	P(D B)
false	false	0.02
false	true	0.98
true	false	0.05
true	true	0.95

A Set of Tables for Each Node

Conditional Probability
Distribution for C given B

B	C	P(C B)
false	false	0.4
false	true	0.6
true	false	0.9
true	true	0.1

For a given combination of values of the parents (B in this example), the entries for $P(C=\text{true} \mid B)$ and $P(C=\text{false} \mid B)$ must add up to 1
eg. $P(C=\text{true} \mid B=\text{false}) + P(C=\text{false} \mid B=\text{false}) = 1$

If you have a Boolean variable with k Boolean parents, this table has 2^{k+1} probabilities (but only 2^k need to be stored)

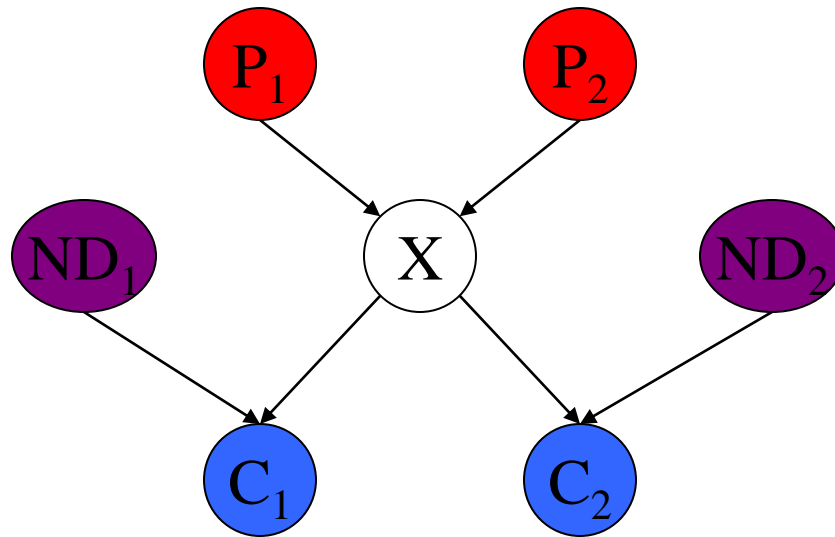
Bayesian Networks

Two important properties:

1. Encodes the conditional independence relationships between the variables in the graph structure
2. Is a compact representation of the joint probability distribution over the variables

Conditional Independence

The Markov condition: given its parents (P_1, P_2), a node (X) is conditionally independent of its non-descendants (ND_1, ND_2)



The Joint Probability Distribution

Due to the Markov condition, we can compute the joint probability distribution over all the variables X_1, \dots, X_n in the Bayesian net using the formula:

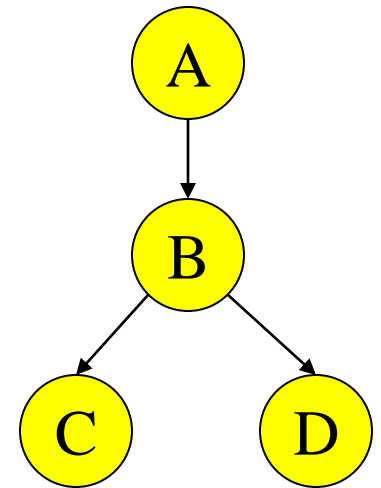
$$P(X_1 = x_1, \dots, X_n = x_n) = \prod_{i=1}^n P(X_i = x_i \mid \text{Parents}(X_i))$$

Where $\text{Parents}(X_i)$ means the values of the Parents of the node X_i with respect to the graph

Using a Bayesian Network Example

Using the network in the example, suppose you want to calculate:

$$\begin{aligned} &P(A = \text{true}, B = \text{true}, C = \text{true}, D = \text{true}) \\ &= P(A = \text{true}) * P(B = \text{true} \mid A = \text{true}) * \\ &\quad P(C = \text{true} \mid B = \text{true}) P(D = \text{true} \mid B = \text{true}) \\ &= (0.4) * (0.3) * (0.1) * (0.95) \end{aligned}$$



Using a Bayesian Network Example

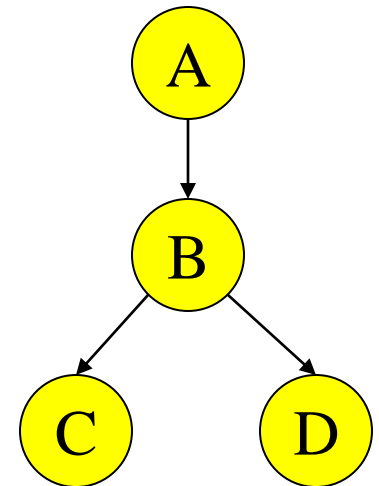
Using the network in the example, suppose you want to calculate:

$$\begin{aligned} &P(A = \text{true}, B = \text{true}, C = \text{true}, D = \text{true}) \\ &= P(A = \text{true}) * P(B = \text{true} \mid A = \text{true}) * \\ &\quad P(C = \text{true} \mid B = \text{true}) P(D = \text{true} \mid B = \text{true}) \\ &= (0.4) * (0.3) * (0.1) * (0.95) \end{aligned}$$

This is from the
graph structure



These numbers are from the
conditional probability tables



Inference

- Using a Bayesian network to compute probabilities is called inference
- In general, inference involves queries of the form:

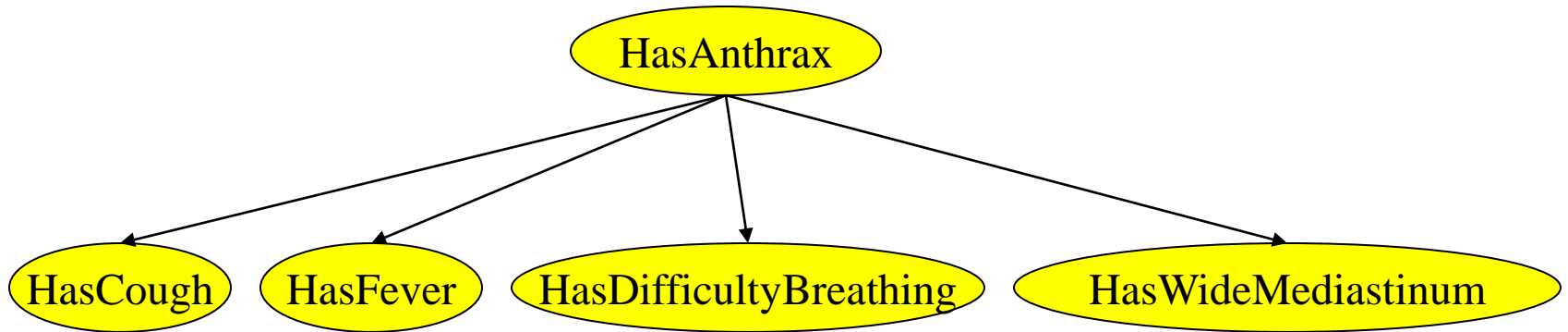
$$P(X \mid E)$$



E = The evidence variable(s)

X = The query variable(s)

Inference



- An example of a query would be:
 $P(\text{HasAnthrax} = \text{true} \mid \text{HasFever} = \text{true}, \text{HasCough} = \text{true})$
- Note: Even though *HasDifficultyBreathing* and *HasWideMediastinum* are in the Bayesian network, they are not given values in the query (ie. they do not appear either as query variables or evidence variables)
- They are treated as unobserved variables

The Bad News

- Exact inference is feasible in small to medium-sized networks
- Exact inference in large networks takes a very long time
- We resort to approximate inference techniques which are much faster and give pretty good results