

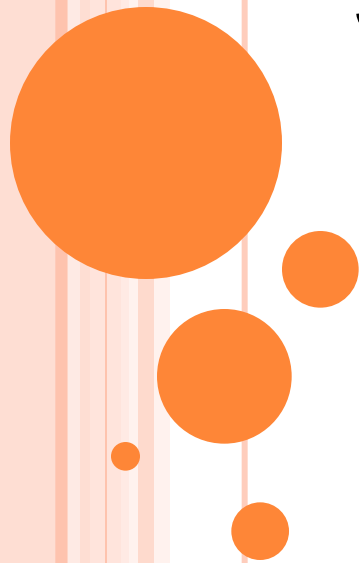


DRONACHARYA
College of Engineering

INTELLIGENT SYSTEMS (CSE-303-F)

Section C

Symbolic reasoning under uncertainty



- Problem are not always consistent, complete and unchanging
- **Monotonicity**
 - New facts can be added to the system. If these new facts are consistent with all the other facts that have already been asserted, then nothing will ever be retracted from the set of facts that are known to be true.



UNCERTAINTY

Several Approaches Related to

- Mathematical and Statistical Theories
- Bayesian Statistics
- Fuzzy Sets



UNCERTAINTY IN AI

Approximate Reasoning, Inexact Reasoning



RELEVANT INFORMATION IS DEFICIENT IN ONE OR MORE

- Information is partial
- Information is not fully reliable
- Representation language is inherently imprecise
- Information comes from multiple sources and it is conflicting
- Information is approximate
- Non-absolute cause-effect relationships exist
- Can include probability in the rules
- IF the interest rate is increasing, THEN the price of stocks will decline (80% probability)



TYPES OF UNCERTAINTY

- Uncertainty in prior knowledge

E.g., some causes of a disease are unknown and are not represented in the background knowledge of a medical-assistant agent



TYPES OF UNCERTAINTY

- Uncertainty in actions

E.g., to deliver this lecture:

I must be able to come to college

my computer must be working

the LCD projector must be working

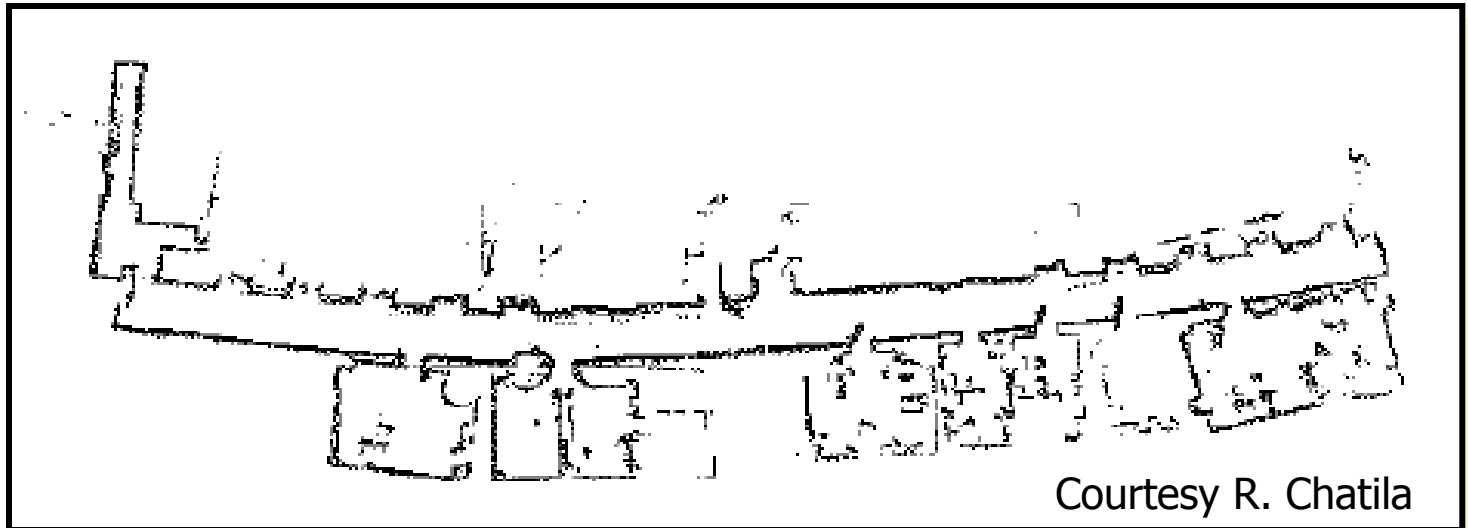
As we discussed with planning, actions are represented with relatively short lists of preconditions, while these lists are in fact arbitrary long. It is not efficient (or even possible) to list all the possibilities.



TYPES OF UNCERTAINTY

- Uncertainty in perception

E.g., sensors do not return exact or complete information about the world; a robot never knows exactly its position.



SOURCES OF UNCERTAINTY

- Laziness (efficiency)
- Ignorance

What we call *uncertainty* is a summary of all that is not explicitly taken into account in the agent's knowledge base (KB).



ASSUMPTIONS OF REASONING WITH PREDICATE LOGIC

(1). Predicate descriptions must be sufficient with respect to the application domain.

Each fact is known to be either true or false. But what does lack of information mean?

- Closed world assumption, assumption based reasoning:

PROLOG: if a fact cannot be proven to be true, assume that it is false

HUMAN: if a fact cannot be proven to be false, assume it is true



ASSUMPTIONS OF REASONING WITH PREDICATE LOGIC (CONT'D)

(2). The information base must be consistent.

- Human reasoning: keep alternative (possibly conflicting) hypotheses. Eliminate as new evidence comes in.



ASSUMPTIONS OF REASONING WITH PREDICATE LOGIC (CONT'D)

(3). Known information grows *monotonically* through the use of inference rules.

○ Need mechanisms to:

- add information based on assumptions (nonmonotonic reasoning), and
- delete inferences based on these assumptions in case later evidence shows that the assumption was incorrect (truth maintenance).



QUESTIONS

- How to represent uncertainty in knowledge?
- How to perform inferences with uncertain knowledge?
- Which action to choose under uncertainty?



APPROACHES TO HANDLING UNCERTAINTY

- Default reasoning [Optimistic]
non-monotonic logic
- Worst-case reasoning [Pessimistic]
adversarial search:
 - *Adversarial search, or game-tree search, is a technique for analyzing an adversarial game in order to try to determine who can win the game and what moves the players should make in order to win.*
- Probabilistic reasoning [Realist]
probability theory



DEFAULT REASONING

- Rationale: The world is fairly normal. Abnormalities are rare.
- So, an agent assumes normality, until there is evidence of the contrary.
- E.g., if an agent sees a bird X, it assumes that X can fly, unless it has evidence that X is a penguin, an ostrich, a dead bird, a bird with broken wings, ...



MODIFYING LOGIC TO SUPPORT NONMONOTONIC INFERENCE

- $p(X) \wedge \text{unless } q(X) \rightarrow r(X)$
- If we
- believe $p(X)$ is true, and
- do not believe $q(X)$ is true (either unknown or believed to be false)
- then we
- can infer $r(X)$
- later if we find out that $q(X)$ is true, $r(X)$ must be retracted

“unless” is a *modal operator*: deals with belief rather than truth



MODIFYING LOGIC TO SUPPORT NONMONOTONIC INFERENCE (CONT'D)

- $p(X) \wedge \text{unless } q(X) \rightarrow r(X)$ in KB
- $p(Z)$ in KB
- $r(W) \rightarrow s(W)$ in KB
- - - - - -
- $\neg q(X) ??$ $q(X)$ is not in KB
- $r(X)$ inferred
- $s(X)$ inferred



EXAMPLE

- If it is snowing and unless there is an exam tomorrow, I can go skiing.
- It is snowing.
- Whenever I go skiing, I stop by at the Chalet to drink hot chocolate.
- - - - - -
- I did not check my calendar but I don't remember an exam scheduled for tomorrow, conclude: I'll go skiing. Then conclude: I'll drink hot chocolate.



“ABNORMALITY”

- $p(X) \wedge \text{unless } ab \ p(X) \rightarrow q(X)$
- ab: abnormal
- Examples: If X is a bird, it will fly unless it is abnormal.
 - (abnormal: broken wing, sick, trapped, ostrich, ...)
- If X is a car, it will run unless it is abnormal.
 - (abnormal: flat tire, broken engine, no gas, ...)



ANOTHER MODAL OPERATOR: M

- $p(X) \wedge M q(X) \rightarrow r(X)$
- If
 - we believe $p(X)$ is true, and
 - $q(X)$ is *consistent with* everything else,
- then we
 - can infer $r(X)$

“M” is a *modal operator* for “is consistent.”



Probabilities and Related Approaches

The Probability Ratio

- $P(X) = \frac{\text{Number of outcomes favoring the occurrence of } X}{\text{Total number of outcomes}}$

Total number of outcomes

- Multiple Probability Values in Many Systems
 - Three-part antecedent (probabilities: 0.9, 0.7, and 0.65)
 - The overall probability:
 $P = (0.9)(0.7)(0.65) = 0.4095$
- Sometimes one rule references another - individual rule probabilities can propagate from one to another

SEVERAL APPROACHES FOR COMBINING PROBABILITIES

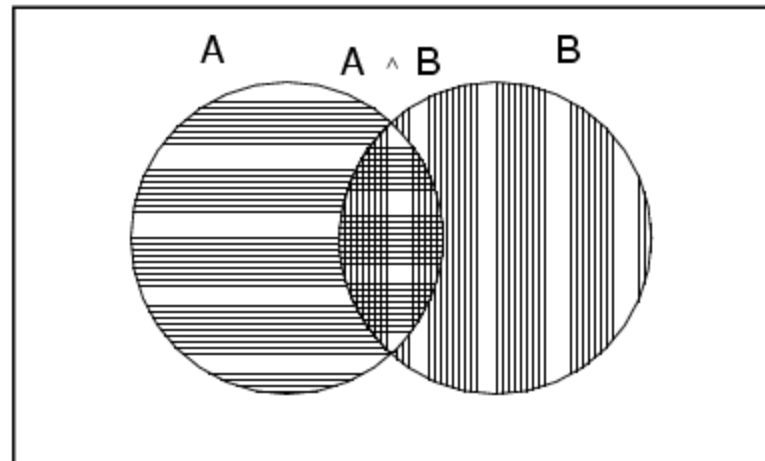
- Probabilities can be
 - **Multiplied (joint probabilities)**
 - **Averaged (simple or a weighted average)**
 - **Highest value**
 - **Lowest value**
- Rules and events are considered independent of each other
- If Dependent - Use the Bayes extension theorem



AXIOMS OF PROBABILITY

- For any propositions A, B
- - $0 \leq P(A) \leq 1$
 - $P(\text{true}) = 1$ and $P(\text{false}) = 0$
 - $P(A \vee B) = P(A) + P(B) - P(A \wedge B)$
 -

True



BAYESIAN APPROACH

- It depend on the use of known prior and likely probabilities.
- It is introduced by Clergyman Thomas Bayes in the eighteenth Century.
- It depends on the use of conditional probabilities of specified events when it is known that other events have occurred.



- For two events H and E with the Probability $P(E) > 0$, the conditional probability of event H, given that event E has occurred, is defined as
 - $P(H/E) = P(H \& E) / P(E)$ 1
- The Conditional probability of event E given that event H occurred can likewise be written
 - $P(E/H) = P(H \& E) / P(H)$ 2
- From 1 & 2
 - $P(H/E) = P(E/H)P(H) / P(E)$ 3



EXAMPLE

- Patient has the certain disease D1 given the symptom E we wish to find out $P(D1/E)$.
 - $P(D1) = 0.05$ and $P(E) = 0.15$ $P(E/D1) = 0.95$
 - $P(D1/E) = 0.32$



- If $P(E)$ is difficult to obtain , then replace H with $\sim H$ in equation 3.

- $P(\sim H/E) = P(E/\sim H)P(\sim H) / P(E)$ 4

- Divide equation $\frac{3}{4}$

- $P(H/E) / P(\sim H/E) = P(E/H)P(H) / P(E/\sim H)P(\sim H)$

- $O(H/E) = L(E/H).O(H)$ 5

- Posterior odds = likelihood ratio * prior odds on H .



SYNTAX

- Basic element: **random variable**
- Similar to propositional logic: possible worlds defined by assignment of values to random variables.
- **Boolean** random variables
 - e.g., *Cavity* (do I have a cavity?)
- **Discrete** random variables
 - e.g., *Weather* is one of *<unny,rainy,cloudy,snow>*
- Domain values must be exhaustive and mutually exclusive

$(A_{25} \text{ gets me there on time} \mid \dots) = 0.04$

$P(A_{90} \text{ gets me there on time} \mid \dots) = 0.70$

$P(A_{120} \text{ gets me there on time} \mid \dots) = 0.95$

$P(A_{1440} \text{ gets me there on time} \mid \dots) = 0.9999$



Suppose I believe the following:

$$P(A_{25} \text{ gets me there on time} \mid \dots) = 0.04$$

$$P(A_{90} \text{ gets me there on time} \mid \dots) = 0.70$$

$$P(A_{120} \text{ gets me there on time} \mid \dots) = 0.95$$

$$P(A_{1440} \text{ gets me there on time} \mid \dots) = 0.9999$$



SYNTAX

- **Atomic event**: A **complete** specification of the state of the world about which the agent is uncertain

E.g., if the world consists of only two Boolean variables *Cavity* and *Toothache*, then there are 4 distinct atomic events:

Cavity = *false* \wedge *Toothache* = *false*

Cavity = *false* \wedge *Toothache* = *true*

Cavity = *true* \wedge *Toothache* = *false*

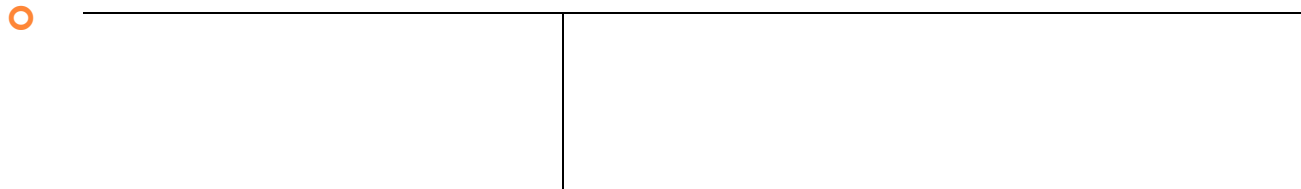
Cavity = *true* \wedge *Toothache* = *true*

Atomic events are mutually exclusive and exhaustive



PRIOR PROBABILITY

- Prior or unconditional probabilities of propositions
- - e.g., $P(\text{Cavity} = \text{true}) = 0.1$ and $P(\text{Weather} = \text{sunny}) = 0.72$ correspond to belief prior to arrival of any (new) evidence
- Probability distribution gives values for all possible assignments:
- - $P(\text{Weather}) = \langle 0.72, 0.1, 0.08, 0.1 \rangle$ (normalized, i.e., sums to 1)
- Joint probability distribution for a set of random variables gives the probability of every atomic event on those random variables



$P(Weather, Cavity)$ = a 4×2 matrix of values:

<i>Weather</i> =	sunny	rainy	cloudy	snow
<i>Cavity</i> = true	0.144	0.02	0.016	0.02
<i>Cavity</i> = false	0.576	0.08	0.064	0.08

- Every question about a domain can be answered by the joint distribution



CONDITIONAL PROBABILITY

- Conditional or posterior probabilities



e.g., $P(\text{cavity} \mid \text{toothache}) = 0.8$

i.e., given that *toothache* is all I know

- (Notation for conditional distributions:



$P(\text{Cavity} \mid \text{Toothache}) = 2\text{-element vector of } 2\text{-element vectors}$)



- If we know more, e.g., *cavity* is also given, then we have
- $$P(\textit{cavity} \mid \textit{toothache}, \textit{cavity}) = 1$$
- New evidence may be irrelevant, allowing simplification, e.g.,
- $P(\textit{cavity} \mid \textit{toothache}, \textit{sunny}) = P(\textit{cavity} \mid \textit{toothache}) = 0.8$
- This kind of inference, sanctioned by domain knowledge, is crucial



CONDITIONAL PROBABILITY

- Definition of conditional probability:
- $P(a \mid b) = P(a \wedge b) / P(b)$ if $P(b) > 0$
- **Product rule** gives an alternative formulation:
 $P(a \wedge b) = P(a \mid b) P(b) = P(b \mid a) P(a)$
- A general version holds for whole distributions, e.g.



$P(\text{Weather}, \text{Cavity}) = P(\text{Weather} \mid \text{Cavity}) P(\text{Cavity})$
(View as a set of 4×2 equations, **not** matrix mult.)

- **Chain rule** is derived by successive application of product rule:
- $$\begin{aligned} P(X_1, \dots, X_n) &= P(X_1, \dots, X_{n-1}) P(X_n \mid X_1, \dots, X_{n-1}) \\ &= P(X_1, \dots, X_{n-2}) P(X_{n-1} \mid X_1, \dots, X_{n-2}) P(X_n \mid X_1, \dots, X_{n-1}) \\ &= \dots \\ &= \prod_{i=1}^n P(X_i \mid X_1, \dots, X_{i-1}) \end{aligned}$$



INFERENCE BY ENUMERATION

- Start with the joint probability distribution:

	<i>toothache</i>		\neg <i>toothache</i>	
	<i>catch</i>	\neg <i>catch</i>	<i>catch</i>	\neg <i>catch</i>
<i>cavity</i>	.108	.012	.072	.008
\neg <i>cavity</i>	.016	.064	.144	.576

- For any proposition ϕ , sum the atomic events where it is true: $P(\phi) = \sum_{\omega: \omega \models \phi} P(\omega)$



INFERENCE BY ENUMERATION

- Start with the joint probability distribution:

	<i>toothache</i>		\neg <i>toothache</i>	
	<i>catch</i>	\neg <i>catch</i>	<i>catch</i>	\neg <i>catch</i>
<i>cavity</i>	.108	.012	.072	.008
\neg <i>cavity</i>	.016	.064	.144	.576

- For any proposition ϕ , sum the atomic events where it is true: $P(\phi) = \sum_{\omega: \omega \models \phi} P(\omega)$
- $P(\text{toothache}) = 0.108 + 0.012 + 0.016 + 0.064 = 0.2$



INFERENCE BY ENUMERATION

- Start with the joint probability distribution:

	<i>toothache</i>		\neg <i>toothache</i>	
	<i>catch</i>	\neg <i>catch</i>	<i>catch</i>	\neg <i>catch</i>
<i>cavity</i>	.108	.012	.072	.008
\neg <i>cavity</i>	.016	.064	.144	.576

- For any proposition ϕ , sum the atomic events where it is true: $P(\phi) = \sum_{\omega: \omega \models \phi} P(\omega)$
- $P(\text{toothache}) = 0.108 + 0.012 + 0.016 + 0.064 = 0.2$



INFERENCE BY ENUMERATION

- Start with the joint probability distribution:

	<i>toothache</i>		\neg <i>toothache</i>	
	<i>catch</i>	\neg <i>catch</i>	<i>catch</i>	\neg <i>catch</i>
<i>cavity</i>	.108	.012	.072	.008
\neg <i>cavity</i>	.016	.064	.144	.576

- Can also compute conditional probabilities:

$$P(\neg \text{cavity} \mid \text{toothache}) = \frac{P(\neg \text{cavity} \wedge \text{toothache})}{P(\text{toothache})}$$

$$= \frac{0.016 + 0.064}{0.108 + 0.012 + 0.016 +$$

0.064

$$= 0.4$$



NORMALIZATION

	<i>toothache</i>		\neg <i>toothache</i>	
	<i>catch</i>	\neg <i>catch</i>	<i>catch</i>	\neg <i>catch</i>
<i>cavity</i>	.108	.012	.072	.008
\neg <i>cavity</i>	.016	.064	.144	.576

- Denominator can be viewed as a **normalization constant** α
- $P(\text{Cavity} \mid \text{toothache}) = \alpha, P(\text{Cavity}, \text{toothache})$
= $\alpha, [P(\text{Cavity}, \text{toothache}, \text{catch}) + P(\text{Cavity}, \text{toothache}, \neg \text{catch})]$
= $\alpha, [<0.108, 0.016> + <0.012, 0.064>]$
= $\alpha, <0.12, 0.08> = <0.6, 0.4>$

General idea: compute distribution on query variable by fixing **evidence variables** and summing over **hidden variables**



INFERENCE BY ENUMERATION, CONTD.

Typically, we are interested in
the posterior joint distribution of the **query variables** Y
given specific values e for the **evidence variables** E

Let the **hidden variables** be $H = X - Y - E$

Then the required summation of joint entries is done by summing out
the hidden variables:

$$P(Y \mid E = e) = \alpha P(Y, E = e) = \alpha \sum_h P(Y, E = e, H = h)$$



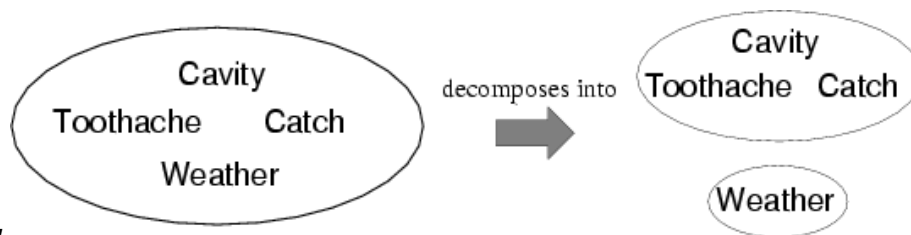
- The terms in the summation are joint entries because Y, E and H together exhaust the set of random variables.
- Obvious problems:
 1. Worst-case time complexity $O(d^n)$ where d is the largest arity
 - 2.
 2. Space complexity $O(d^n)$ to store the joint distribution
 - 3.
 3. How to find the numbers for $O(d^n)$ entries?



INDEPENDENCE

- A and B are independent iff

$$P(A/B) = P(A) \quad \text{or} \quad P(B/A) = P(B) \quad \text{or} \quad P(A, B) = P(A) P(B)$$



$$P(\textit{Toothache}, \\ = P(\textit{Toothache}, \textit{Catch}, \textit{Cavity}) P(\textit{Weather})$$

- 32 entries reduced to 12; for n independent biased coins, $O(2^n) \rightarrow O(n)$
- Absolute independence powerful but rare
- Dentistry is a large field with hundreds of variables, none of which are independent. What to do?



CONDITIONAL INDEPENDENCE

- $P(\textit{Toothache}, \textit{Cavity}, \textit{Catch})$ has $2^3 - 1 = 7$ independent entries
- If I have a cavity, the probability that the probe catches in it doesn't depend on whether I have a toothache:
- (1) $P(\textit{catch} \mid \textit{toothache}, \textit{cavity}) = P(\textit{catch} \mid \textit{cavity})$
- The same independence holds if I haven't got a cavity:
- (2) $P(\textit{catch} \mid \textit{toothache}, \neg \textit{cavity}) = P(\textit{catch} \mid \neg \textit{cavity})$
- *Catch* is **conditionally independent** of *Toothache* given *Cavity*:
- $P(\textit{Catch} \mid \textit{Toothache}, \textit{Cavity}) = P(\textit{Catch} \mid \textit{Cavity})$



CONDITIONAL INDEPENDENCE CONTD.

- Write out full joint distribution using chain rule:

- $P(\textit{Toothache}, \textit{Catch}, \textit{Cavity})$
 $= P(\textit{Toothache} \mid \textit{Catch}, \textit{Cavity}) P(\textit{Catch}, \textit{Cavity})$

 $= P(\textit{Toothache} \mid \textit{Catch}, \textit{Cavity}) P(\textit{Catch} \mid \textit{Cavity})$
 $P(\textit{Cavity})$

 $= P(\textit{Toothache} \mid \textit{Cavity}) P(\textit{Catch} \mid \textit{Cavity}) P(\textit{Cavity})$

I.e., $2 + 2 + 1 = 5$ independent numbers

In most cases, the use of conditional independence reduces the size of the representation of the joint distribution from exponential in n to linear in n .

Conditional independence is our most basic and robust form of knowledge about uncertain environments.



- Equivalent statements:

$$P(\textit{Toothache} \mid \textit{Catch}, \textit{Cavity}) = P(\textit{Toothache} \mid \textit{Cavity})$$

$$P(\textit{Toothache}, \textit{Catch} \mid \textit{Cavity}) = P(\textit{Toothache} \mid \textit{Cavity}) P(\textit{Catch} \mid \textit{Cavity})$$



BAYES' RULE

- Product rule $P(a \wedge b) = P(a | b) P(b) = P(b | a) P(a)$
 \Rightarrow Bayes' rule: $P(a | b) = P(b | a) P(a) / P(b)$
- or in distribution form
- $P(Y|X) = P(X|Y) P(Y) / P(X) = \alpha P(X|Y) P(Y)$



- Useful for assessing **diagnostic** probability from **causal** probability:
 - $P(\text{Cause}|\text{Effect}) = P(\text{Effect}|\text{Cause}) P(\text{Cause}) / P(\text{Effect})$
 - E.g., let M be meningitis, S be stiff neck:
 - $P(m|s) = P(s|m) P(m) / P(s) = 0.8 \times 0.0001 / 0.1 = 0.0008$
 - Note: posterior probability of meningitis still very small!



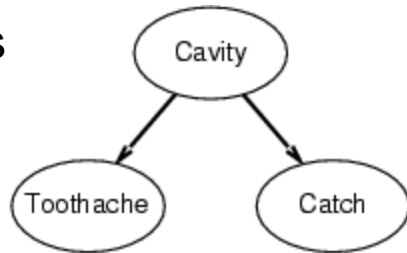
BAYES' RULE AND CONDITIONAL INDEPENDENCE

$$\begin{aligned} P(\text{Cavity} \mid \text{toothache} \wedge \text{catch}) \\ &= \alpha P(\text{toothache} \wedge \text{catch} \mid \text{Cavity}) P(\text{Cavity}) \\ &= \alpha P(\text{toothache} \mid \text{Cavity}) P(\text{catch} \mid \text{Cavity}) P(\text{Cavity}) \end{aligned}$$

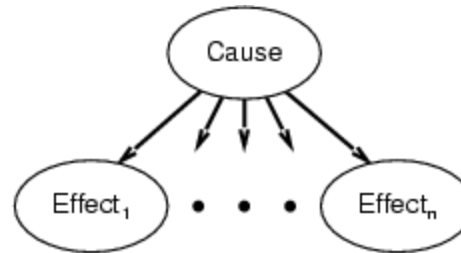
- This is an example of a **naïve Bayes** model:

○

$P(\text{Caus}$



$se)$



- Total number of parameters is **linear** in n

SUMMARY

- Probability is a rigorous formalism for uncertain knowledge
- Joint probability distribution specifies probability of every atomic event
- Queries can be answered by summing over atomic events
- For nontrivial domains, we must find a way to reduce the joint size
- Independence and conditional independence provide the tools



DEMPSTER-SHAFER THEORY

○ Drawbacks of using Bayesian theory

- The probabilities are described as a single numeric point value.
 - Distortion to precision that is actually available for supporting evidence.
 - When we assert with probability 0.7 that the dollar will fall against the Japanese Yen over the next six months, what we really mean is we have a fairly strong conviction there is a chance of about 0.6 to 0.8 say, that it will fall



DEMPSTER-SHAFER THEORY

○ Drawbacks of using Bayesian theory

- No way to differentiate between ignorance and uncertainty.
- Example
 - One of the three A, B, C terrorist group has planted a bomb. Let C found guilty and $P(C) = 0.8$. According to traditional theory rest of the probability will distribute amongst other without having any knowledge about them.



DEMPSTER-SHAFER THEORY

○ Drawbacks of using Bayesian theory

- Forced to regard belief and disbelief as functional opposite.
 - Ex. If $P(A) = 0.3$ then $P(\sim A) = 0.7$ so that $P(A) + P(\sim A) = 1$

AS A REMEDY FOR THE ABOVE PROBLEMS,
GENERALISES THEORY HAS BEEN PROPOSED
BY ARTHUR DEMPSTER (1968) AND EXTENDED
BY STUDENT GLENN SHAFER (1976).



DEMPSTER-SHAFER THEORY

- Separate probability masses may be assigned to all subsets of a universe of discourse rather than just to individual single members.
 - It Permit the inequality $P(A) + P(\sim A) \leq 1$
 - It assume a universe of discourse U
 - A set corresponding to n proposition, exactly one of which is true.
 - The Propositions are assumed to be exhaustive and mutually exclusive.



DEMPSTER-SHAFER THEORY

- Let 2^U denote all subsets of U .
- Let the set function m defined on 2^U , be a mapping to $[0,1]$,
 - $m : 2^U \rightarrow [0,1]$, be, such that for all subsets $A \subset U$
 - $m(\emptyset) = 0$
 - $\sum_{A \subset U} m(A) = 1$
- The function m defines a probability distribution on 2^U . (not just on the singletons of U as in classical theory)
- It represent the measure of belief committed exactly to A .
- IT IS POSSIBLE TO ASSIGN BELIEF TO EACH SUBSET A OF U WITHOUT ASSIGNING ANY TO ANYTHING SMALLER.



DEMPSTER-SHAFER THEORY

- Bel(A) = $\sum_{B \subset A} m(B)$

A BELIEF FUNCTION, BEL , CORRESPONDING TO A SPECIFIC m FOR THE SET A, IS DEFINED AS THE SUM OF BELIEFS COMMITTED TO EVERY SUBSET OF A BY m.

Ex: if U contain the mutually exclusive subsets A, B, C and D then

$$\begin{aligned} \text{Bel}(\{A, C, D\}) = & m(\{A, C, D\}) + m(\{A, C\}) + m(\{A, D\}) + m(\{C, D\}) \\ & + m(\{A\}) + m(\{C\}) + m(\{D\}) \end{aligned}$$



DEMPSTER-SHAFER THEORY

○ Some related terms and facts:

- In D-S T, a belief interval can also be defined for a subset A . It is represented as the subinterval $[\text{Bel}(A), \text{Pl}(A)]$ of $[0,1]$.
- **Support:** $\text{Bel}(A)$ is also called support of A :
- **Plausibility:** $\text{Pl}(A) = 1 - \text{Bel}(\sim A)$, the plausibility of A .
- **Focal Element:** The subsets A of U are called the focal elements of the support function Bel when $m(A) > 0$.



DEMPSTER-SHAFER THEORY

○ Some related terms and facts:

- **Doubts:** $\text{Bel}(A)$ partially describes the beliefs about proposition A , belief in $\sim A$ (doubt) can be defined as $D(A) = \text{Bel}(\sim A)$
- Hence Upper bound of the interval can be defined as $P1(A) = 1 - D(A) = 1 - \text{Bel}(\sim A)$.
- **Confidence:** The Belief interval, $[\text{Bel}(A), P1(A)]$, is also sometimes referred to as a confidence in A , while the quantity $P1(A) - \text{Bel}(A)$ is referred to as the **uncertainty in A** .



DEMPSTER-SHAFER THEORY

- It can be shown that (Prade 1983)
 - $P1(\$) = 0, P1(U) = 1$
 - For all A
 - $P1(A) \geq Bel(A),$
 - $Bel(A) + Bel(\sim A) \leq 1,$
 - $P1(A) + P1(\sim A) \geq 1,$ and
 - For ACB,
 - $Bel(A) \leq Bel(B), P1(A) \leq P1(B)$
- In interpreting the above definitions, it should be noted that a portion of belief may be committed to a set of propositions, but need not be, and if committed, it is not necessary to its negation. However, a belief committed to a proposition is committed to any other proposition it implies.



DEMPSTER-SHAFER THEORY

- $[0,1]$ represents no belief in support of the proposition
- $[0,0]$ represents the belief the proposition is false
- $[1,1]$ represents the belief the proposition is true
- $[\cdot 3,1]$ represents partial belief in the proposition
- $[0,\cdot 8]$ represents partial disbelief in the proposition
- $[\cdot 2,\cdot 7]$ represents belief from evidence both for and against the proposition



DEMPSTER-SHAFER THEORY

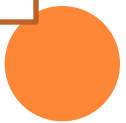
- When evidence is available from two or more independent KS Bel1 and Bel2 then Demster has provided such a combining function denoted by $\text{Bel1} \circ \text{Bel2}$.
 - Let m_1 and m_2 are two basic probability assignment function to the Bel1 and Bel2
 - Let A_1, \dots, A_k be focal elements for Bel1 and B_1, \dots, B_p be the focal elements for Bel2.
 - $m_1(A_i)$ and $m_2(B_j)$ each assign probability masses on the unit interval. They can be orthogonally represented as:



DEMPSTER-SHAFER THEORY

$m_1(A_1) \dots \dots \dots m_1(A_i) \dots \dots$

.				
.				
$m_2(B_j)$				
.				
.				
$m_2(B_1)$				



DEMPSTER-SHAFER THEORY

- The unit square represents the total probability mass assigned by both m_1 and m_2 for all of their common subsets.
- Particular rectangle within the square, shown as the intersection of the sets A_i and B_j , has committed to it the measure $m_1(A_i)m_2(B_j)$. Therefore the total probability mass committed to C will be

- $\sum_{A_i \wedge B_j = C} m_1(A_i)m_2(B_j)$ Its normalised form after removing

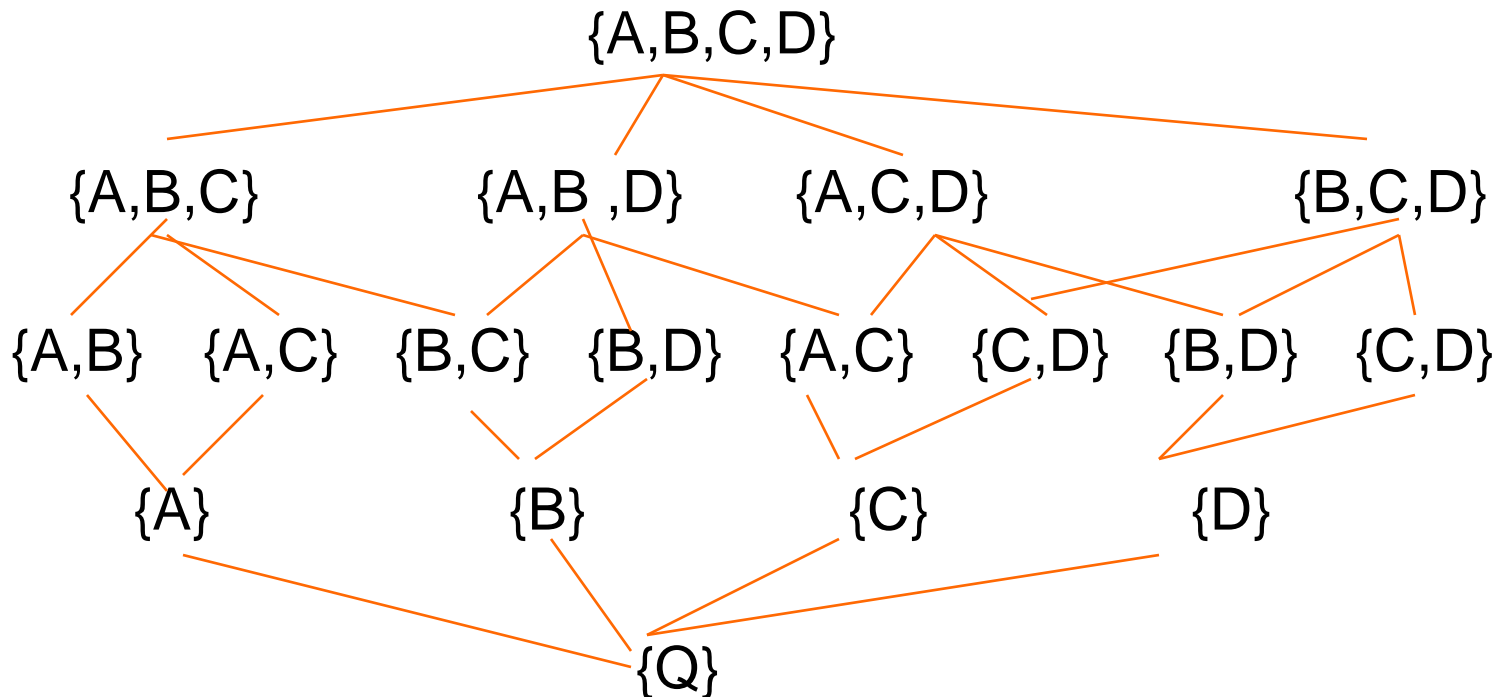
$$A_i \wedge B_j = @$$

$$m_1 \circ m_2 = \sum_{A_i \wedge B_j} m_1(A_i)m_2(B_j) / \sum_{A_i \wedge B_j \neq \emptyset} m_1(A_i)m_2(B_j)$$



DEMPSTER-SHAFER THEORY

- Ex of problem of identifying the terrorist group from A, B, C, D THE POSSIBLE CASES ARE



DEMPSTER-SHAFER THEORY

- Assume that one group A and C were responsible to a degree of $m_1(\{A,C\}) = 0.6$,
- ANOTHER SOURCE OF EVIDENCE DISPROVES THE BELIEF THAT c WAS INVOLVED hence it means $m_2(\{A,B,D\}) = 0.7$.
- To obtain the pooled evidence, we compute the following quantities.
 - $m_1 \circ m_2(\{A\}) = 0.6 * 0.7 = 0.42$
 - $m_1 \circ m_2(\{A,C\}) = 0.6 * 0.3 = 0.18$
 - $m_1 \circ m_2(\{A,B,D\}) = 0.4 * 0.7 = 0.28$
 - $m_1 \circ m_2(\{U\}) = 0.4 * 0.3 = 0.12$
 - $m_1 \circ m_2$ for all other set = 0
 - $Bel_1(\{A,C\}) = m(\{A,C\}) + m(\{A\}) + m(\{C\})$



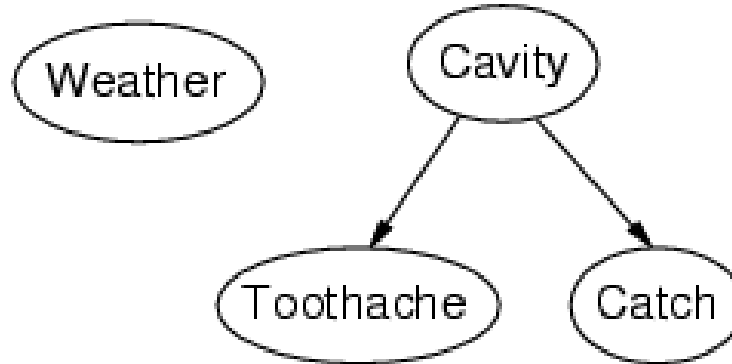
BAYESIAN NETWORKS

- A simple, graphical notation for conditional independence assertions and hence for compact specification of full joint distributions
- Syntax:
 - a set of nodes, one per variable
 -
 - a directed, acyclic graph (link \approx "directly influences")
 - a conditional distribution for each node given its parents:
$$\mathbf{P}(X_i \mid \text{Parents}(X_i))$$
- In the simplest case, conditional distribution represented as a **conditional probability table** (CPT) giving the distribution over X_i for each combination of parent values



EXAMPLE

- Topology of network encodes conditional independence assertions:



- *Weather* is independent of the other variables
- *Toothache* and *Catch* are conditionally independent given *Cavity*

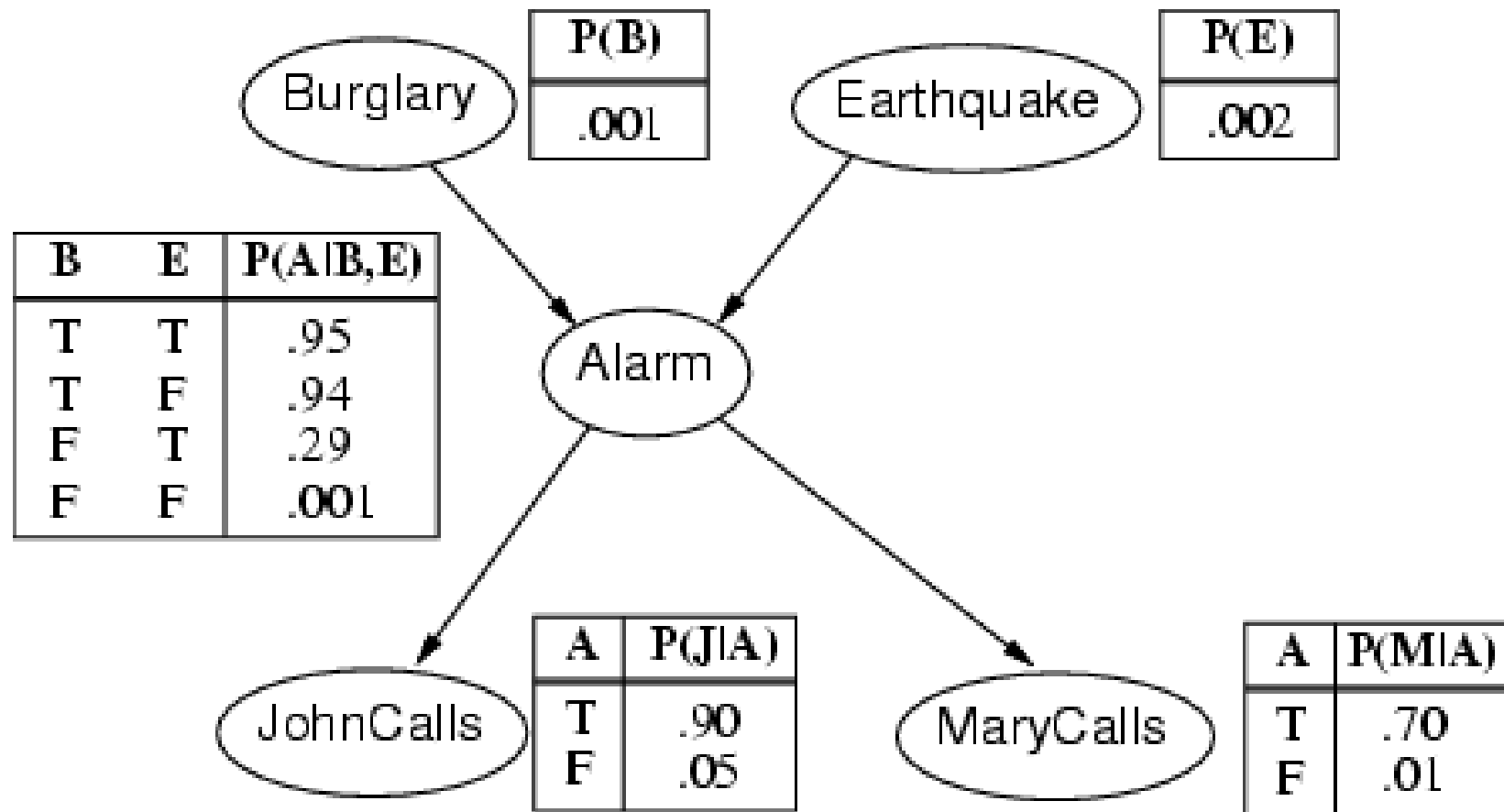


EXAMPLE

- I'm at work, neighbor John calls to say my alarm is ringing, but neighbor Mary doesn't call. Sometimes it's set off by minor earthquakes. Is there a burglar?
- Variables: *Burglary*, *Earthquake*, *Alarm*, *JohnCalls*, *MaryCalls*
- Network topology reflects "causal" knowledge:
 - A burglar can set the alarm off
 - An earthquake can set the alarm off
 - The alarm can cause Mary to call
 - The alarm can cause John to call

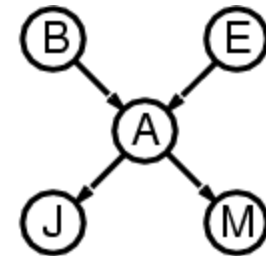


EXAMPLE CONTD.



COMPACTNESS

- A CPT for Boolean X_i with k Boolean parents has 2^k rows for the combinations of parent values
- Each row requires one number p for $X_i = \text{true}$ (the number for $X_i = \text{false}$ is just $1-p$)
- If each variable has no more than k parents, the complete network requires $O(n \cdot 2^k)$ numbers
- I.e., grows linearly with n , vs. $O(2^n)$ for the full joint distribution
- For burglary net, $1 + 1 + 4 + 2 + 2 = 10$ numbers (vs. $2^5 - 1 = 31$)

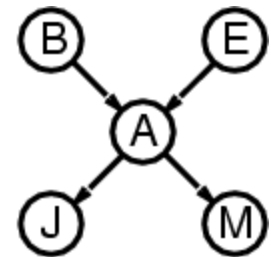


SEMANTICS

The full joint distribution is defined as the product of the local conditional distributions:

n

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



e.g., $P(j \wedge m \wedge a \wedge \neg b \wedge \neg e)$

$$= P(j \mid a) P(m \mid a) P(a \mid \neg b, \neg e) P(\neg b) P(\neg e)$$



CONSTRUCTING BAYESIAN NETWORKS

- 1. Choose an ordering of variables X_1, \dots, X_n
- 2. For $i = 1$ to n
 - add X_i to the network

-

- select parents from X_1, \dots, X_{i-1} such that

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

This choice of parents guarantees:

$$P(X_1, \dots, X_n) = \prod_{i=1}^n P(X_i | X_1, \dots, X_{i-1})$$

(chain rule)

$$= \prod_{i=1}^n P(X_i | \text{Parents}(X_i))$$

(by construction)



EXAMPLE

- Suppose we choose the ordering M, J, A, B, E



MaryCalls

JohnCalls

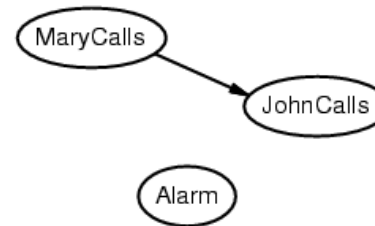
$$P(J \mid M) = P(J)?$$



EXAMPLE

- Suppose we choose the ordering M, J, A, B, E

-



$$P(J \mid M) = P(J)?$$

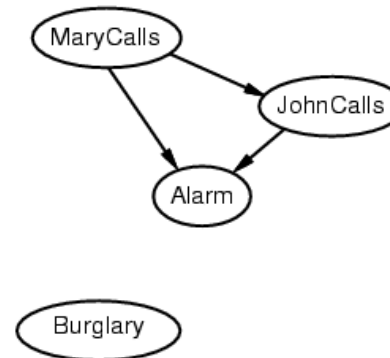
No

$$P(A \mid J, M) = P(A \mid J)? \quad P(A \mid J, M) = P(A)?$$



EXAMPLE

- Suppose we choose the ordering M, J, A, B, E



$$P(J \mid M) = P(J)?$$

No

$$P(A \mid J, M) = P(A \mid J)? \quad P(A \mid J, M) = P(A)? \quad \text{No}$$

$$P(B \mid A, J, M) = P(B \mid A)?$$

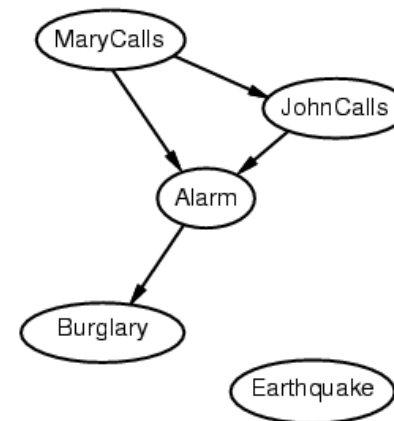
$$P(B \mid A, J, M) = P(B)?$$



EXAMPLE

- Suppose we choose the ordering M, J, A, B, E

-



$$P(J \mid M) = P(J)?$$

No

$$P(A \mid J, M) = P(A \mid J)? \quad P(A \mid J, M) = P(A)? \quad \text{No}$$

$$P(B \mid A, J, M) = P(B \mid A)? \quad \text{Yes}$$

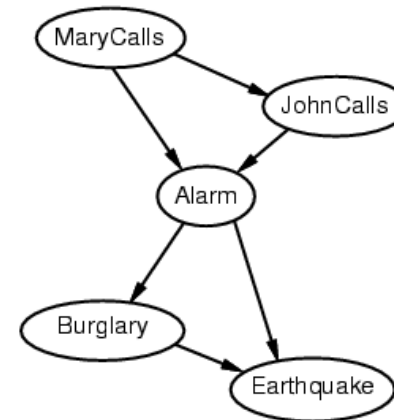
$$P(B \mid A, J, M) = P(B)? \quad \text{No}$$

$$P(E \mid B, A, J, M) = P(E \mid A)?$$

EXAMPLE

- Suppose we choose the ordering M, J, A, B, E

-



$$P(J \mid M) = P(J)?$$

No

$$P(A \mid J, M) = P(A \mid J)? \quad P(A \mid J, M) = P(A)? \quad \text{No}$$

$$P(B \mid A, J, M) = P(B \mid A)? \quad \text{Yes}$$



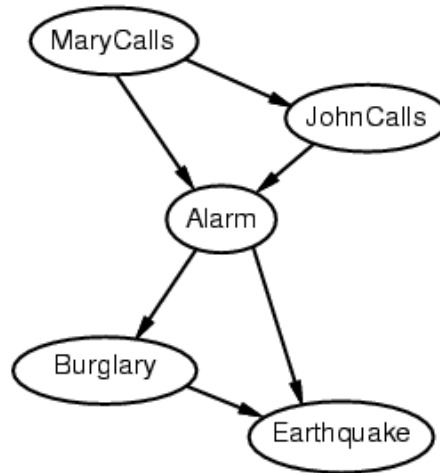
$P(B \mid A, J, M) = P(B)$? No

$P(E \mid B, A, J, M) = P(E \mid A)$? No

$P(E \mid B, A, J, M) = P(E \mid A, B)$? Yes



EXAMPLE CONTD.



- Deciding conditional independence is hard in noncausal directions
- (Causal models and conditional independence seem hardwired for humans!)
- Network is less compact: $1 + 2 + 4 + 2 + 4 = 13$ numbers needed
-

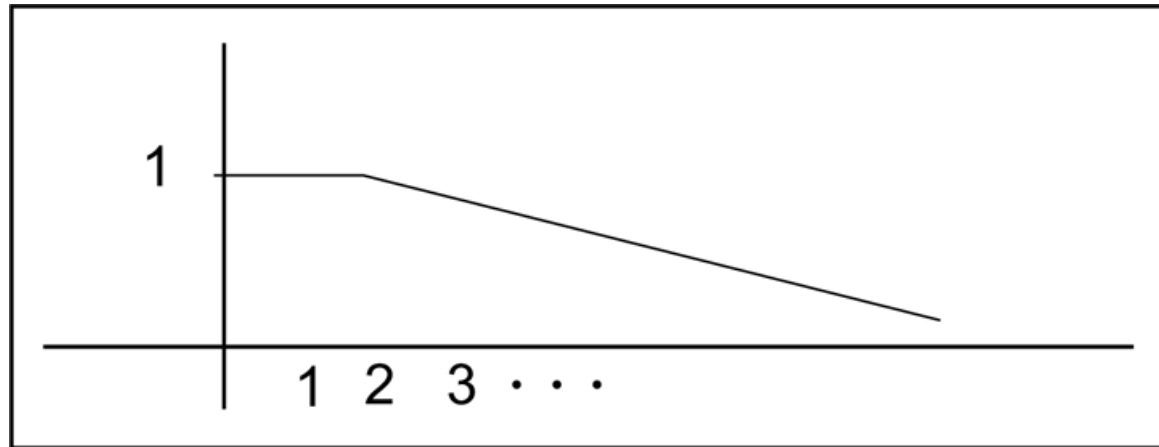


SUMMARY

- Bayesian networks provide a natural representation for (causally induced) conditional independence
- Topology + CPTs = compact representation of joint distribution
- Generally easy for domain experts to construct



THE FUZZY SET REPRESENTATION FOR “SMALL INTEGERS”

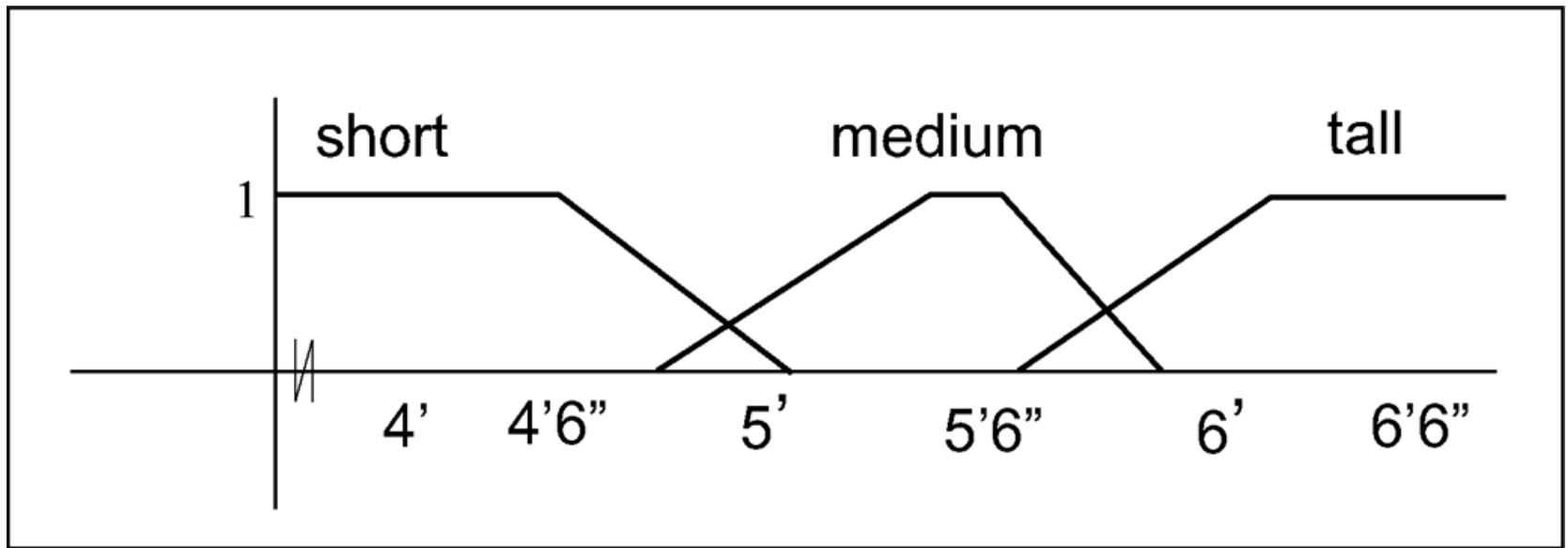


REASONING WITH FUZZY SETS

- Lotfi Zadeh's fuzzy set theory
- Violates two basic assumption of set theory
 - For a set S , an element of the universe either belongs to S or the complement of S .
 - For a set S , and element cannot belong to S or the complement S at the same time
- Jack is 5'7". Is he tall? Does he belong to the set of tall people? Does he not belong to the set of tall people?



A FUZZY SET REPRESENTATION FOR THE SETS SHORT, MEDIAN, AND TALL MALES

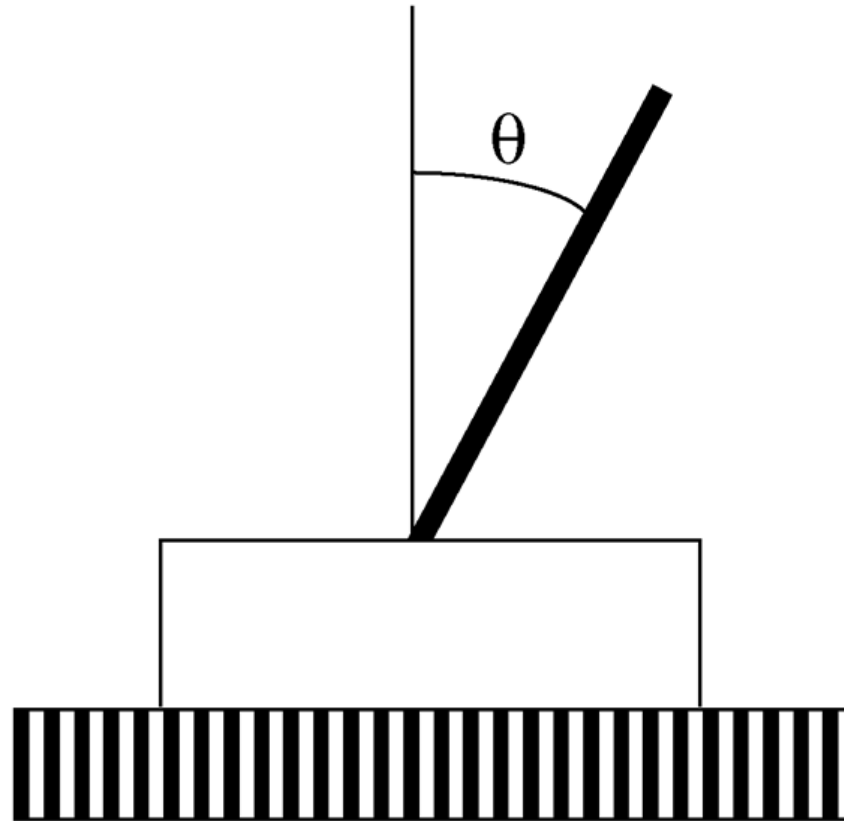


FUZZY LOGIC

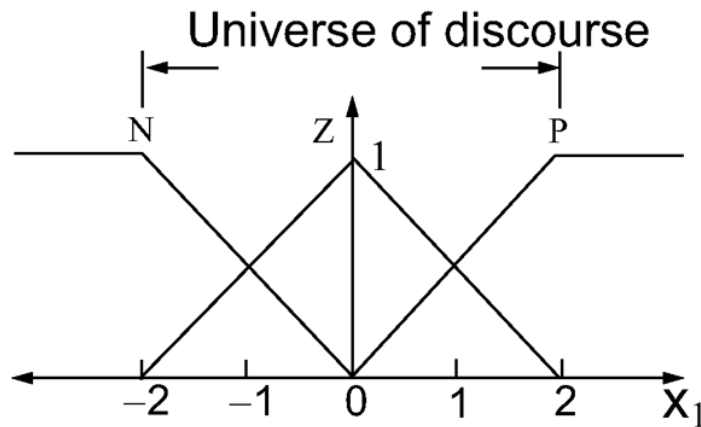
- Provides rules about evaluating a fuzzy truth, T
- The rules are:
 - $T(A \wedge B) = \min(T(A), T(B))$
 - $T(A \vee B) = \max(T(A), T(B))$
 - $T(\neg A) = 1 - T(A)$
- Note that unlike logic $T(A \vee \neg A) \neq T(\text{True})$



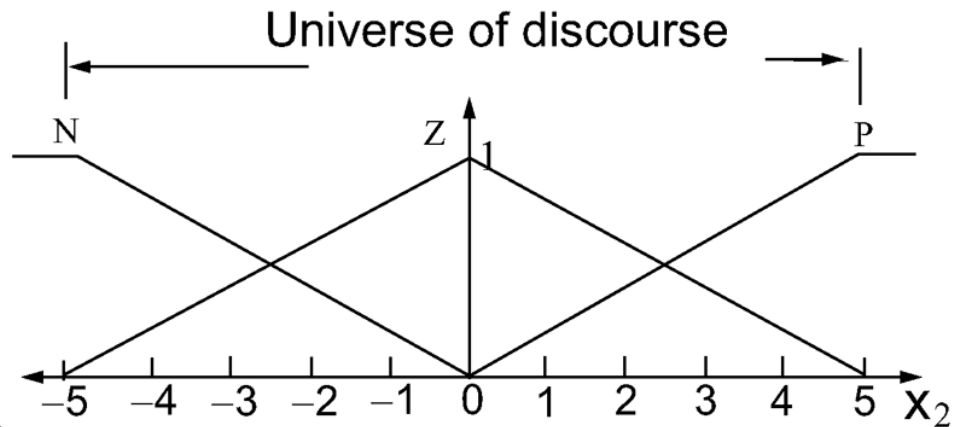
THE INVERTED PENDULUM AND THE ANGLE θ AND $D\theta/DT$ INPUT VALUES.



THE FUZZY REGIONS FOR THE INPUT VALUES θ (A) AND $D\theta/DT$ (B)



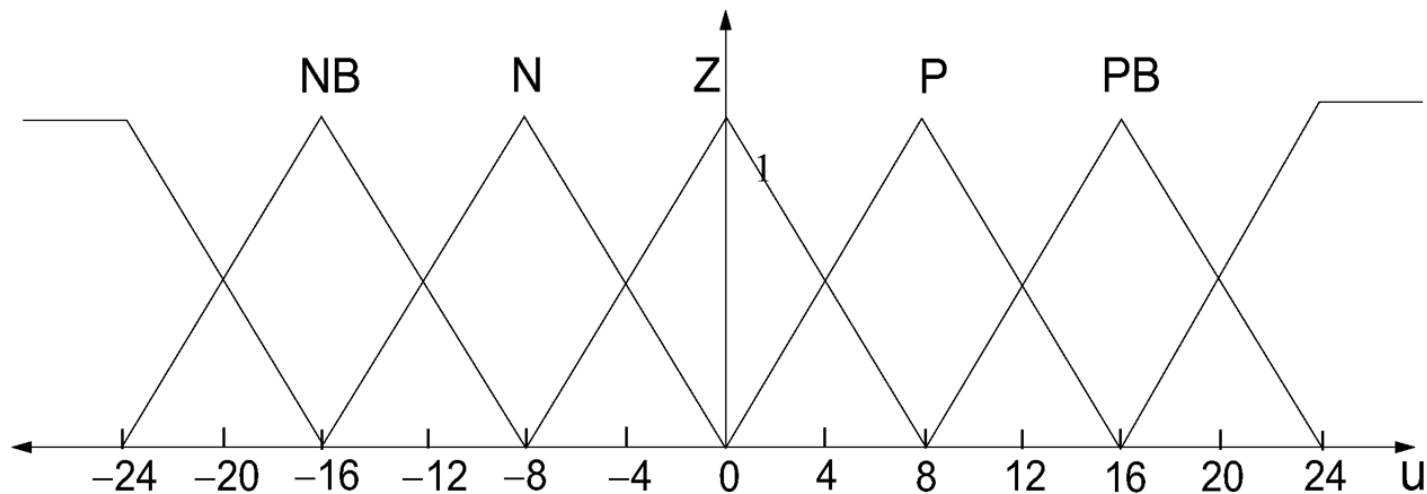
a.



b.

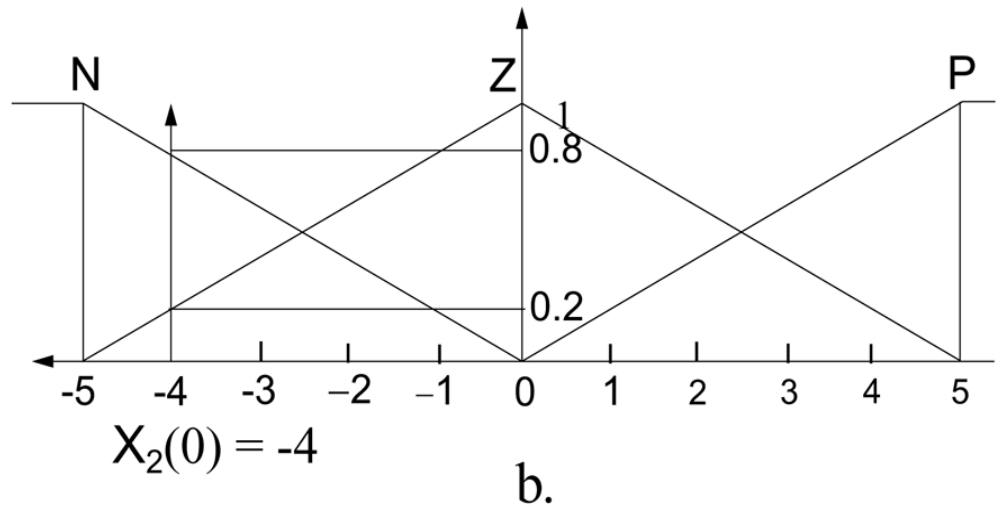
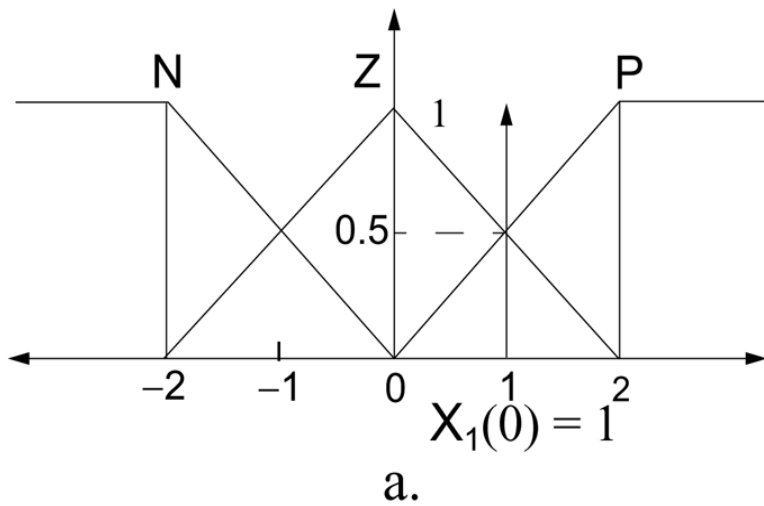


THE FUZZY REGIONS OF THE OUTPUT VALUE u , INDICATING THE MOVEMENT OF THE PENDULUM BASE



THE FUZZIFICATION OF THE INPUT MEASURES

$x_1=1, x_2 = -4$

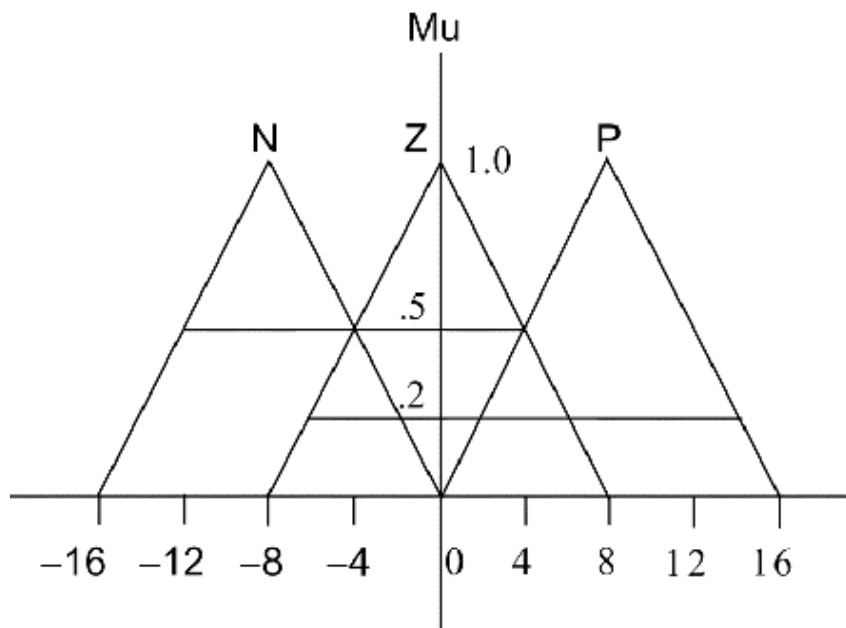


THE FUZZY ASSOCIATIVE MATRIX (FAM) FOR THE PENDULUM PROBLEM

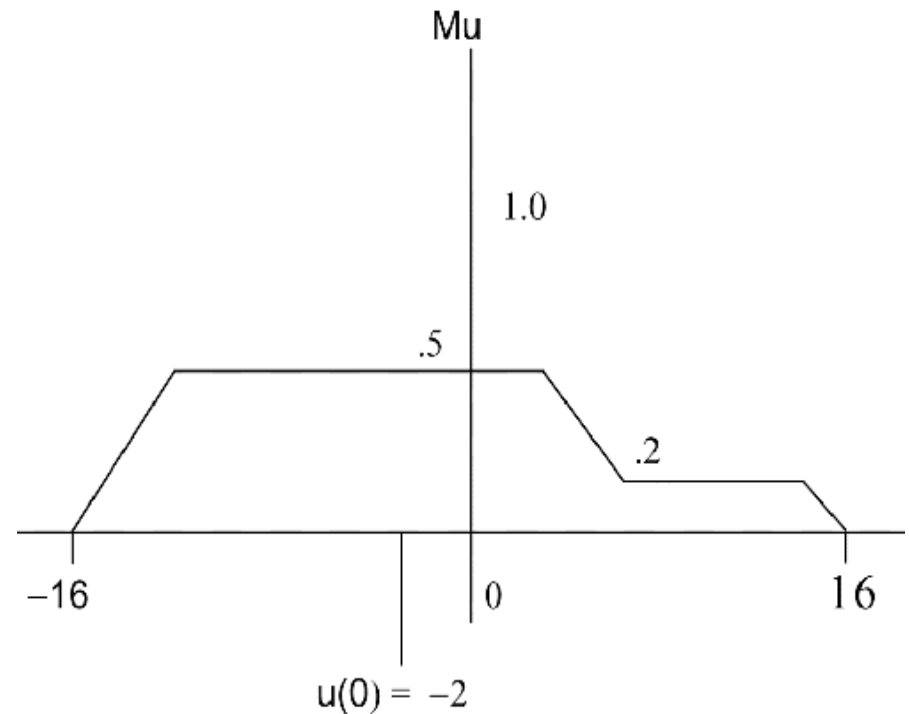
$x_1 \backslash x_2$	P	Z	N
P	PB	P	Z
Z	P	Z	N
N	Z	N	NB



THE FUZZY CONSEQUENTS (A), AND THEIR UNION (B)



a.



b.

The centroid of the union (-2) is the crisp output.



MINIMUM OF THEIR MEASURES IS TAKEN AS
THE MEASURE OF THE RULE RESULT

IF $x_1 = P$ AND $x_2 = Z$ THEN $u = P$
 $\min(0.5, 0.2) = 0.2 P$

IF $x_1 = P$ AND $x_2 = N$ THEN $u = Z$
 $\min(0.5, 0.8) = 0.5 Z$

IF $x_1 = Z$ AND $x_2 = Z$ THEN $u = Z$
 $\min(0.5, 0.2) = 0.2 Z$

IF $x_1 = Z$ AND $x_2 = N$ THEN $u = N$
 $\min(0.5, 0.8) = 0.5 N$



PROCEDURE FOR CONTROL

- Take the crisp output and fuzzify it
- Check the Fuzzy Associative Matrix (FAM) to see which rules fire
(4 rules fire in the example)
- Find the rule results
 - ANDed premises: take minimum
 - ORed premises: take maximum
- Combine the rule results
(union in the example)
- Defuzzify to obtain the crisp output
(centroid in the example)



COMMENTS

- “fuzzy” refers to sets (as opposed to *crisp* sets)
- Fuzzy logic is useful in engineering control where the measurements are imprecise
- It has been successful in commercial control applications: automatic transmissions, trains, video cameras, electric shavers
- useful when there are small rule bases, no chaining of inferences, tunable parameters
- The theory is not concerned about how the rules are created, but how they are combined
- The rules are not chained together, instead all fire and the results are combined

