

Uncertainty and Probabilistic Reasoning

Overview

- Uncertainty
- Decision theory example
- Probability basics
- Conditional probability
- Axioms of probability
- Joint probability distribution
- Bayes rule
- Bayes rule: Example

1

2

Uncertainty

- Problem with first-order logic: agents almost never have full access to the whole truth about their environment.
- Therefore, the agent must act under **uncertainty**.
- Uncertainty can also arise because of incompleteness and incorrectness in the agent's understanding of the properties in the environment.
- Incomplete, because there are too many conditions to explicitly enumerate.

There are trade-offs (playing safe can result in other annoyances), thus the right thing to do depends on both the relative importance of various goals and the likelihood (and degree to which) they will be achieved.

3

Example: Trying to Catch a Flight

A_t : plan to leave home t minutes before the flight departure time.

- The traveler needs to make a decision in an uncertain environment: car can break down, traffic can be extremely congested, natural disaster, etc.
- Such worst-case scenarios are hard to explicitly enumerate: the list goes on – ran out of gas, spouse/children in an emergency, flight crews goes on a strike, etc. etc.
- Thus the traveler only has an incomplete understanding of the situation.
- The traveler can play safe by going with plan A_{1440} , but this can cause the traveler to wait for a long time at the airport before departure.

4

Difficulties in Applying F-O-L in Uncertain Domains

For example, application of first-order logic in medical diagnosis domain can fail because of these reasons:

- Laziness: cannot list the complete set of antecedents and consequents needed to ensure an exceptionless rule, and too hard to use the enormous rules that result.
- Theoretical ignorance: medical science has no complete theory.
- Practical ignorance: even though we have all the rules, it is practically impossible to run all the tests.

Similar situation arises in law, business, dating, etc. The agent's knowledge can at best provide only a **degree of belief**. **Probability theory** is well suited for such a domain.

5

Example

When playing black jack,

- as new cards are drawn and shown, your degree of belief in the fact that you need more cards can change.

What about poker? or slot machine?

7

Acquisition of New Information and Probability

- The degree of belief changes as an agent perceives or acquires new information from the world: we call this the **evidence**.
- This is analogous to saying whether or not a given logical sentence is entailed by (i.e. is a logical consequence of) the knowledge base, because the truth value can change when new facts are added to the KB.
- Before the evidence is received, we talk about **prior** or **unconditional** probability.
- After the evidence is obtained, we talk about **posterior** or **conditional** probability.

6

Rational Decisions Under Uncertainty: Decision Theory

- There are trade-offs, and an agent must first have **preferences** between different results when a certain plan was executed.
- **Utility theory** deals with such preferences: how useful is such and such result to the agent?
- **Decision theory** is a general theory of rational decision under uncertainty, combining **probability theory** and **utility theory**.

8

Decision Theory

- An agent is rational iff it chooses the action that yields the highest expected utility, averaged over all possible outcomes of the action:

Principle of Maximum Expected Utility

- Example: backgammon (discussed earlier) – min-max trees with probabilistic levels.

9

Decision Theory: Example

Decision theory = Probability theory + Utility theory

	Utility of Resulting State	Probability
Action 1	10	0.2
Action 2	10000	0.001
Action 3	5	0.799

Which action would an optimal Decision Theoretic Agent take?

11

Decision Theoretic Agent

function DT-Agent (*percept*) **returns** *action*

static: a set probabilistic *belief* about the state of the world

calculate updated probabilities for current state based on *percept* and past actions

calculate outcome probabilities for actions, given action descriptions and prob of current states.

select *action* with highest expected utility given prob of outcomes and utility information.

return *action*

10

Decision Theory: Example

Decision theory = Probability theory + Utility theory

	Utility of Resulting State \times Probability	Expected Utility
Action 1	10×0.2	2
Action 2	1000×0.001	1
Action 3	5×0.799	3.995 \leftarrow

Action 3 has the maximum expected utility, thus action 3 will be carried out.

12

Probability: Notations^a

- **Random variable**: variable that can take on different values
 - A, B, \dots : boolean values (**T** or **F**).
 - X, Y, \dots : numerical values or other multi-valued enumerations (1, 2, 0.5, Cloudy, Rainy, Sunny, ...)
- $P(X = v)$: **probability** of the variable X having value v .
 - This can be viewed as an *event*.
 - For boolean variables, $P(A)$ means $P(A = \mathbf{T})$, and $P(\neg A)$ means $P(A = \mathbf{F})$.
- $\mathbf{P}(X)$: **probability distribution**, a full list of probabilities for all possible values that X can take (note that **P** is in **bold**).

^aAll conventions follow Russel & Norvig

13

Examples

- Boolean:
 $P(\text{Infected}) = 0.01, P(\neg \text{Infected}) = 0.99$.
- Multi valued:
 $P(\text{Dice} = 1) = \frac{1}{6}, P(\text{Dice} = 2) = \frac{1}{6}, \dots$
- Multi valued:
 $P(\text{Weather} = \text{Sunny}) = 0.7,$
 $P(\text{Weather} = \text{Rainy}) = 0.2, \dots$

14

Logical Connectives and Conditional Probability

- Logical connectives can be used:
 $P(A \vee B), P(A \wedge \neg B), P(\text{Cavity} \wedge \neg \text{Insured}), \text{etc.}$
- Conditional Probability $P(A|B)$ (read *probability of A given B*):

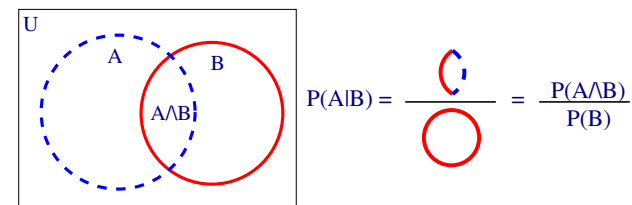
$$P(\text{Cavity}|\text{Toothache}) = 0.8$$

- As new evidence comes in, the conditional probability gets updated:

$$P(\text{Cavity}|\underbrace{\text{Toothache} \wedge \text{BadBreath}})$$

15

Conditional Probability



- Think about the **area** occupied by each event.
- The bounding rectangle U has an area of 1, thus

$$P(A) = \frac{\text{Area of } A}{\text{Area of } U} = \frac{\text{Area of } A}{1} = \text{Area of } A$$

- $P(A|B)$ means B now takes on the role of U . Within this limited event space, what is the probability of A .

16

The Axioms of Probability

All axioms

1. All probabilities are between 0 and 1

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

2. For a valid proposition A (**T** under all interpretations):

$P(A) = 1$, and for an inconsistent proposition A (**F** under all interpretations): $P(A) = 0$.

3. $P(A \vee B) = P(A) + P(B) - P(A \wedge B)$

Other properties follow from these three axioms.

17

Joint Probability Distribution

For random variables X_1, X_2, \dots, X_n ,

- An **atomic event** is an assignment of particular values to each random variable.
- The **joint probability distribution** $\mathbf{P}(X_1, X_2, \dots, X_n)$ completely specifies the probabilities of all **atomic events**.
- Thus,

$$\left[\sum_{(v_1, v_2, \dots, v_n) \in \mathbf{V}} P(X_1 = v_1, X_2 = v_2, \dots, X_n = v_n) \right] = 1,$$

where \mathbf{V} is a set of all possible n -vectors that the vector (X_1, X_2, \dots, X_n) can assume..

19

Other Properties

- From the axioms,

$$P(A \vee \neg A) = P(A) + P(\neg A) - P(A \wedge \neg A)$$

$$P(\mathbf{T}) = P(A) + P(\neg A) - P(\mathbf{F})$$

$$1 = P(A) + P(\neg A)$$

$$P(\neg A) = 1 - P(A)$$

- More generally, the **sum** of probabilities $P(X = v)$ is 1, for all values v the random variable X can take:

$$\left[\sum_{v \in V} P(X = v) \right] = 1,$$

where V is the set of all possible values X can take.

18

Joint Probability Distribution: Example

	Toothache	\neg Toothache	Sum
Cavity	0.04	0.06	$P(C) = 0.1$
\neg Cavity	0.01	0.89	$P(\neg C) = 0.9$
Sum	$P(T) = 0.05$	$P(\neg T) = 0.95$	$\sum = 1.0$

Abbreviations: $C = \text{Cavity}, T = \text{Toothache}$

- $P(C \vee T) = P(C) + P(T) - P(C \wedge T) = 0.1 + 0.05 - 0.04 = 0.11$
- $P(C|T) = \frac{P(C \wedge T)}{P(T)} = \frac{0.04}{0.05} = 0.8$
- $P(T|C) = \frac{P(C \wedge T)}{P(C)} = \frac{0.04}{0.1} = 0.5$

In practice, writing a full joint probability table like this is impossible (or too much effort): for n boolean random variables, you need 2^n entries.

20

Bayes' Rule

- From $P(A|B) = \frac{P(A \wedge B)}{P(B)}$ and $P(B|A) = \frac{P(A \wedge B)}{P(A)}$, we get

$$P(A|B)P(B) = P(B|A)P(A)$$

and in turn from which we get the **Bayes' Rule**:

$$P(B|A) = \frac{P(A|B)P(B)}{P(A)}$$

21

Example: Application of Bayes' Rule

Exercise 14.3 (p. 433): After your yearly checkup, the doctor has bad news and good news. The bad news is that you tested positive for a serious disease, and that the test is 99% accurate (i.e. the probability of testing positive given that you have the disease is 0.99, as is the probability of testing negative given you don't have the disease). The good news is that this is a rare disease, striking only 1 in 10,000 people. Why is it good news that the disease is rare? What are the chances that you actually have the disease?

T : tested positive, $\neg T$: tested negative, D : have disease, $\neg D$: clean.

23

Extended Bayes' Rule

$$\mathbf{P}(Y|X, E) = \frac{\mathbf{P}(X|Y, E)\mathbf{P}(Y|E)}{\mathbf{P}(X|E)}$$

This rule follows from

$$\mathbf{P}(A, B|E) = \mathbf{P}(A|B, E)\mathbf{P}(B|E):$$

$$\begin{aligned}\mathbf{P}(A, B|E) &= \frac{\mathbf{P}(A, B, E)}{\mathbf{P}(E)} \\ &= \frac{\mathbf{P}(A|B, E)\mathbf{P}(B, E)}{\mathbf{P}(E)} \\ &= \mathbf{P}(A|B, E)\mathbf{P}(B|E)\end{aligned}$$

Note: $\mathbf{P}(Y|X, E) = \mathbf{P}(Y|\underbrace{(X, E)})$.

Exercise: text book, exercise 14.5b and 14.6 (p. 434).

22

Solution: Good News and Bad News

These are given:

$$\begin{aligned}P(T|D) &= 0.99 \\ P(\neg T|\neg D) &= 0.99 \\ P(D) &= \frac{1}{10,000} = 0.0001\end{aligned}$$

We want to calculate the probability that you have the disease given a positive test result:

$$P(D|T)$$

We can use Bayes' rule to derive this probability.

24

Solution: Good News and Bad News (cont'd)

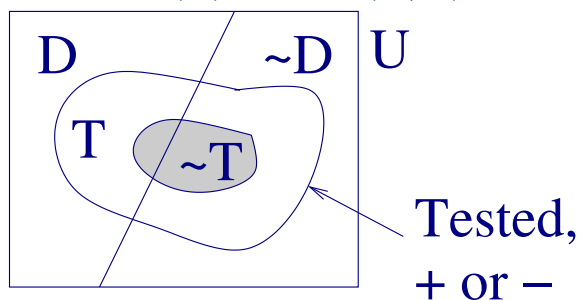
$$P(D|T) = \frac{P(T|D)P(D)}{P(T)}$$

- $P(T|D) = 0.99$, $P(\neg T|\neg D)$, and $P(D) = 0.0001$ are given.
- From these, we can get $P(\neg T|D) = 0.01$, $P(T|\neg D) = 0.01$, and $P(\neg D) = 0.9999$.

Since $P(T|D)$ and $P(D)$ are given, we only need to calculate $P(T)$.

25

Calculating $P(T)$ given $P(T|D)$ and $P(D)$



$$\begin{aligned} P(T) &= P(T \wedge D) + P(T \wedge \neg D) \\ &= P(T|D)P(D) + P(T|\neg D)P(\neg D) \end{aligned}$$

- $\{D\} \cup \{\neg D\}$ completely account for the whole population, but $\{T\} \cup \{\neg T\}$ does not cover the whole population (because you **did not test everyone!**).

27

Solution: Good News and Bad News (cont'd)

Observation ^a: $P(T) = P(T|D)P(D) + P(T|\neg D)P(\neg D)$.

Thus, $P(T) = 0.99 \times 0.0001 + 0.01 \times 0.9999 = 0.010098$, and with this,

$$P(D|T) = \frac{0.99 \times 0.0001}{0.010098} = 0.0098,$$

which is slightly less than 1%.

Exercise: how accurate should the test be so that $P(D|T)$ is greater than 0.95 (i.e. 95%)?

^a $P(T) = P(T \wedge D) + P(T \wedge \neg D)$.

26

Calculating $P(T)$ given $P(T|D)$ and $P(D)$

Another way of deriving

$$P(T) = P(T|D)P(D) + P(T|\neg D)P(\neg D):$$

From

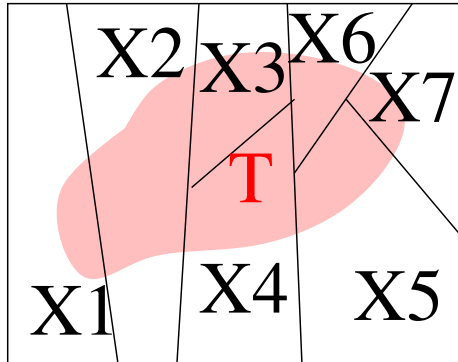
$$\begin{aligned} P(D|T) &= \frac{P(T|D)P(D)}{P(T)} \\ P(\neg D|T) &= \frac{P(T|\neg D)P(\neg D)}{P(T)} \end{aligned}$$

and from $P(D|T) + P(\neg D|T) = 1$,

$$\begin{aligned} 1 &= \frac{P(T|D)P(D)}{P(T)} + \frac{P(T|\neg D)P(\neg D)}{P(T)}, \text{ thus} \\ P(T) &= P(T|D)P(D) + P(T|\neg D)P(\neg D) \end{aligned}$$

28

Calculating $P(T)$: General Case



More Generally, if $\left[\sum_{x \in \{x_1, x_2, \dots, x_n\}} P(X = x) \right] = 1$ and events $X = x_m$ and $X = x_n$ are disjoint for all $m \neq n$,

$$P(T) = \sum_{x \in \{x_1, x_2, \dots, x_n\}} P(T | \underbrace{X = x}) P(X = x)$$

29

Key Points

- Application of theorem proving: question answering
- Uncertainty
- Decision theory example: how prob theory and decision theory are combined
- Probability basics: terminology, notations.
- Joint probability distribution: concept
- Conditional probability: definition, various ways of representing conditional prob.
- Axioms of probability: basic axioms, and using them to prove simple equalities.
- Bayes rule: definition and application.

What's The Big Deal?

- $P(T|D)$ may be easier to obtain: you can run the test on a small pool of known patients (say 100) at a hospital.
- $P(D|T)$ is much harder to obtain directly. Since the test makes 1 mistake out of 100 tests, if you run the test on 10,000 people, you'll get 100 false-positives, and one genuine patient who tests positive (consider that $P(T) = 0.010098$). So, just to get about 100 people testing positive, you have to run the tests on 10,000 people.
- $P(D)$ serves as a **prior** in this case. In many cases, the prior represents subjective **belief** of the person calculating the probability in case $P(D)$ is not directly measurable.

30

Overview

- Diagnostic vs. causal knowledge
- Calculating $P(T)$ given $P(T|D)$ and $P(D)$
- Ratios of conditional probabilities and causes of an phenomenon
- Example: object recognition
- Bayesian updating

Diagnostic vs. Causal Knowledge

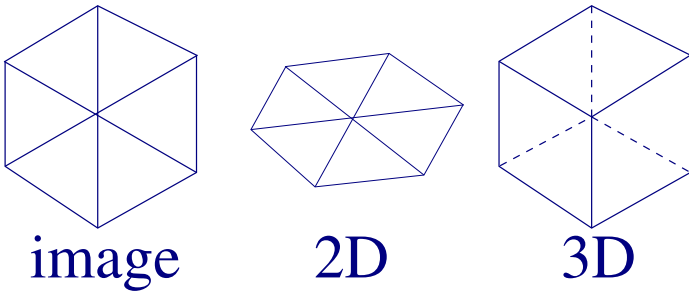
Consider these probabilities:

- $P(\text{Symptom}|\text{Disease})$: **causal knowledge**
- relatively fixed.
- $P(\text{Disease})$: somewhat variable.
- $P(\text{Disease}|\text{Symptom})$: **diagnostic knowledge**
- fluctuates as $P(\text{Disease})$ change.

$P(\text{Disease}|\text{Symptom})$ directly measured can be no longer accurate when $P(\text{Disease})$ changes (e.g. an epidemic outburst), however the calculation based on Bayes' rule can be much more robust.

33

Example: The Problem of Object Recognition



Given an image projected on the retina, what is the more likely cause? the 2D hexagon? or a transparent 3D cube? This is basically a computer vision problem.

$$\frac{P(\text{Hexagon}|\text{Image})}{P(\text{Cube}|\text{Image})} = \frac{P(\text{Image}|\text{Hexagon})P(\text{Hexagon})}{P(\text{Image}|\text{Cube})P(\text{Cube})} = \frac{a}{b}$$

A probabilistic vision agent can make a decision based on such a ratio.

35

Comparison of Conditional Probabilities

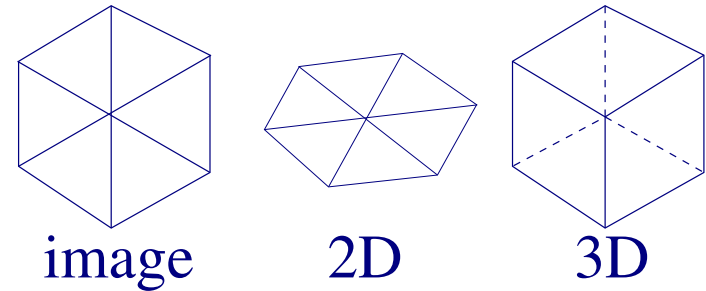
When C_1 or C_2 can cause phenomenon (or effect) E , to find out the which is the more probable cause of phenomenon E , we do not need to explicitly calculate $P(E)$:

- $P(C_1|E) = \frac{P(E|C_1)P(C_1)}{P(E)}$
- $P(C_2|E) = \frac{P(E|C_2)P(C_2)}{P(E)}$
- From the above, we get:

$$\frac{P(C_1|E)}{P(C_2|E)} = \frac{P(E|C_1)P(C_1)}{P(E|C_2)P(C_2)} = \frac{a}{b}$$

34

Example: Object Recognition (cont'd)

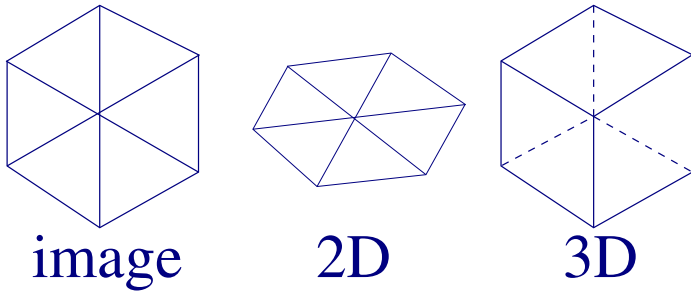


$$\frac{P(\text{Hexagon}|\text{Image})}{P(\text{Cube}|\text{Image})} = \frac{P(\text{Image}|\text{Hexagon})P(\text{Hexagon})}{P(\text{Image}|\text{Cube})P(\text{Cube})} = \frac{a}{b}$$

- Decision: if $a/b > 1$, it is most likely that a hexagon generated the image. If $a/b < 1$, it is most likely that a cube generated the image.

36

Example: Object Recognition (cont'd)



$$\frac{P(\text{Hexagon}|\text{Image})}{P(\text{Cube}|\text{Image})} = \frac{P(\text{Image}|\text{Hexagon})P(\text{Hexagon})}{P(\text{Image}|\text{Cube})P(\text{Cube})} = \frac{a}{b}$$

- Why is $P(\text{Image}|\text{Hexagon})$ easier to calculate than $P(\text{Hexagon}|\text{Image})$?
- What about $P(\text{Hexagon})$ and $P(\text{Cube})$?

37

Bayesian Updating

An Alternative: gradually work in the multiple evidences – **Bayesian Updating**

- Reformulate the Bayes' rule so that conditional probability of events given combined evidences (such as $P(A|B \wedge C)$) are not necessary.
- Use domain knowledge to replace the more complex conditional probabilities with known, simpler ones (utilize **conditional independence**).

Bayesian updating makes combining evidences efficient (more detail next time).

39

Combining Multiple Evidences

Suppose we have these conditional probabilities

$P(\text{Cavity}|\text{Toothache})$ and $P(\text{Cavity}|\text{Catch})$. What if we want to know $P(\text{Cavity}|\text{Toothache} \wedge \text{Catch})$? These are the alternatives:

- Look up the joint probability table: not practical or even impossible in most cases
- We can calculate

$$P(\text{Cav}|\text{Ache} \wedge \text{Catch}) = \frac{P(\text{Ache} \wedge \text{Catch}|\text{Cav})P(\text{Cav})}{P(\text{Ache} \wedge \text{Catch})}$$

but, calculating the new conditional prob and the normalization factor is a pain.

38

Bayesian Updating: Example

We want to calculate $P(\text{Cavity}|\text{Ache} \wedge \text{Catch})$:

$$\begin{aligned} P(\text{Cav}|\text{Ache} \wedge \text{Catch}) &= P(\text{Cav}|\text{Ache}) \frac{P(\text{Catch}|\text{Ache} \wedge \text{Cav})}{P(\text{Catch}|\text{Ache})} \\ &= \underbrace{P(\text{Cav}) \frac{P(\text{Ache}|\text{Cav})}{P(\text{Ache})}}_{\text{old}} \underbrace{\frac{P(\text{Catch}|\text{Ache} \wedge \text{Cav})}{P(\text{Catch}|\text{Ache})}}_{\text{new}} \end{aligned}$$

Problem is that $P(\text{Catch}|\text{Ache} \wedge \text{Cav})$ may be equally hard to calculate. However, we can make these assumptions (**Conditional Independence of Ache and Catch given Cav**):

$$\begin{aligned} P(\text{Catch}|\text{Cav} \wedge \text{Ache}) &= P(\text{Catch}|\text{Cav}) \\ P(\text{Ache}|\text{Cav} \wedge \text{Catch}) &= P(\text{Ache}|\text{Cav}) \end{aligned}$$

Only thing that remains is

$P(\text{Ache})P(\text{Catch}|\text{Ache}) = P(\text{Catch} \wedge \text{Ache})$, which can be eliminated by normalization (**Exercise:** try this – see Exercise 14.7).

Bayesian Updating: Example (cont'd)

So, after replacing the factors:

$$P(Cav|Ache \wedge Catch) = \alpha \underbrace{P(Cav)P(Ache|Cav)}_{old} \underbrace{P(Catch|Cav)}_{new},$$

where α is the normalization constant needed to make $P(Cav|Ache \wedge Catch)$ add up to 1.

41

Key Points

- How is subjective belief utilized in Bayesian analysis?
- Bayesian updating: why does it make probabilistic inference efficient when multiple evidence comes in?

43

Bayesian Updating: Summary

X and Y are independent given Z (conditional independence):

$$P(X|Y, Z) = P(X|Z)$$

Simplified Bayes' rule for multiple evidence is^a:

$$P(Z|X, Y) = \alpha P(Z)P(X|Z)P(Y|Z),$$

where α is the normalization constant.

Thus, Bayesian Updating makes combining multiple evidence easy.

^a **Note:** Z is the cause, and X and Y are the effects.

42

Overview

- Belief network

44

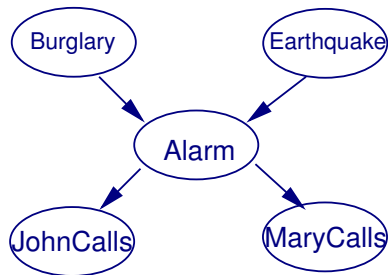
Probabilistic Reasoning

Belief Network represent the dependence between random variables, and give a concise specification of the joint probability distribution. It is represented as a directed acyclic graph (DAG):

1. a set of random variables : nodes of the network
2. a set of directed edges from one node to another
3. each node has a conditional probability table that quantifies the effect the parents have on that node. The parents are the nodes pointing to that node.
4. the graph has no cycles

45

Belief Network: Example (cont'd)

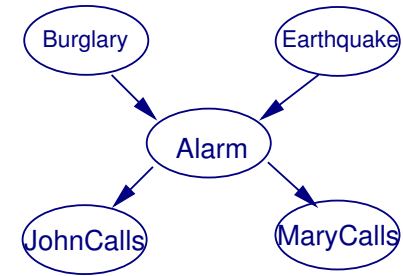


Example question: If you got calls from John and Mary, what is the chance of it being a totally false alarm (not a burglary, nor an earthquake)?

You can ask any conjunctive combination.

47

Belief Network: Example

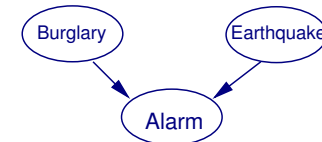


New burglar alarm was installed.

- The alarm can be triggered by either an actual burglary or an earthquake.
- Neighbors John and Mary agreed to call you at work when they hear the alarm.

46

Belief Network: Example (cont'd)



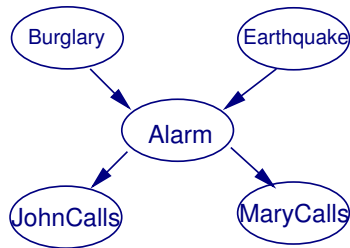
Each node has a conditional probability table fully describing $\mathbf{P}(Current|Parent_1, Parent_2, \dots, Parent_n)$:

Burglary	Earthquake	P	¬P
T	T	0.950	0.050
T	F	0.950	0.050
F	T	0.290	0.710
F	F	0.001	0.999

$\mathbf{P} = \mathbf{P}(Alarm|Burglary, Earthquake)$

48

Semantics of Belief Networks

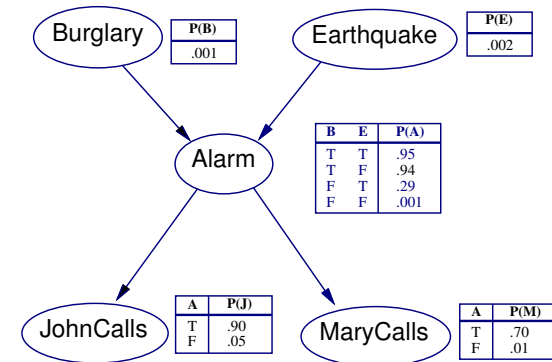


The network can be viewed as

- a representation of the **joint probability distribution** (this view is helpful when **constructing** the network), or
- an encoding of a collection of **conditional independence** statements (this view is helpful when designing effective **inference procedures**).

49

Belief Network: Representing Joint Prob. Dist.



Each row in the conditional probability tables are:

$$\begin{aligned}
 P(X_1 = x_1, \dots, X_n = x_n) &= P(x_1, x_2, \dots, x_n) \\
 &= \prod_{i=1}^n P(x_i | Parent(X_i))
 \end{aligned}$$

50

Belief Network: Representing Joint Prob. Dist.

(cont'd)

$$\begin{aligned}
 P(X_1 = x_1, \dots, X_n = x_n) &= P(x_1, x_2, \dots, x_n) \\
 &= \prod_{i=1}^n P(x_i | Parent(X_i))
 \end{aligned}$$

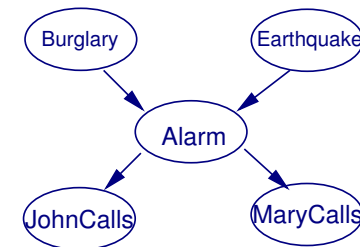
$Parent(X_j)$ refers to the event when $X_j = x_j$.

Imagine a case where each node in the example has a **T** or **F** assignment. X_j will then be either **T** or **F** for all j .

The belief network fully defines a joint probability distribution!

51

Calculating Probability of a Joint Event

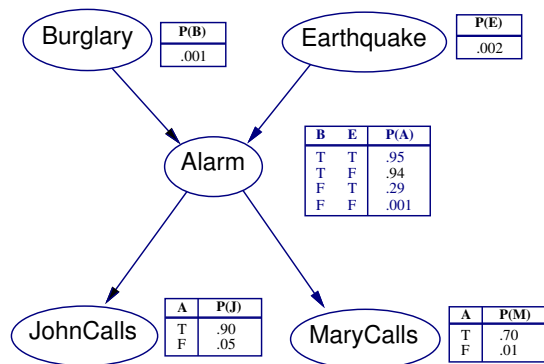


Calculate the probability of the event that the alarm (A) has sounded but neither a burglary ($\neg B$) nor an earthquake ($\neg E$) occurred, and both John (J) and Mary (M) call:

$$\begin{aligned}
 &P(J \wedge M \wedge A \wedge \neg B \wedge \neg E) \\
 &= P(J | Prnts(J))P(M | Prnts(M))P(A | Prnts(A))P(\neg B)P(\neg E) \\
 &= P(J | A)P(M | A)P(A | \neg B \wedge \neg E)P(\neg B)P(\neg E)
 \end{aligned}$$

52

Calculating Probability of a Joint Event (cont'd)



$$P(J|A)P(M|A)P(A|\neg B \wedge \neg E)P(\neg B)P(\neg E)$$

$$= 0.9 \times 0.7 \times 0.001 \times 0.999 \times 0.998 = 0.00062$$

53

Overview

- Constructing a belief network
- Inference in belief networks
- Knowledge engineering

55

Key Points

- Belief network: definition, semantics, extracting probabilities of certain conjunction of events.

54

Joint Probability Distribution Under Conditional Independence

In Belief Networks, the joint probability is given as follows:

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

This is derived from the two following equations:

$$\begin{aligned}
 P(x_1, \dots, x_n) &= P(x_n | x_{n-1}, \dots, x_1) P(x_{n-1}, \dots, x_1) \\
 &= P(x_n | x_{n-1}, \dots, x_1) P(x_{n-1} | x_{n-2}, \dots, x_1) \\
 &\quad \dots P(x_2 | x_1) P(x_1) \\
 &= \prod_{i=1}^n P(x_i | x_{i-1}, \dots, x_1) \tag{1}
 \end{aligned}$$

$$P(X_i | X_{i-1}, \dots, X_1) = P(X_i | \text{Parents}(X_i)) \tag{2}$$

56

Belief Network Construction

Given the nodes, we have to find which nodes are **directly** influenced by certain nodes, and from this, find out **conditional independence**.

- For example:

$$\mathbf{P}(MaryCalls|JohnCalls, Alarm, Earthquake, Burglary) = \mathbf{P}(MaryCalls|Alarm)$$
- So, even if *Earthquake* and *Burglary* can be found up-stream (e.g. *Earthquake* causing *Alarm* to go off, in turn causing *MaryCalls*), those events are conditionally independent from *MaryCalls*.

To construct a Belief Network, we need to find such dependency structure.

57

Evaluation of the Construction Algorithm

- Because newly added nodes cannot point to existing nodes, the resulting graph is always **acyclic**.
- Violation of axioms of probability is avoided.
- Compact, compared to the full joint probability table (**locally structured, or sparse**).
 - Belief network with n nodes (binary variables) and k parents per node has $n2^k$ entries in the conditional probability tables:

$$n \text{ nodes} \times 2^k \text{ per each node.}$$

- Full joint probability table: 2^n

59

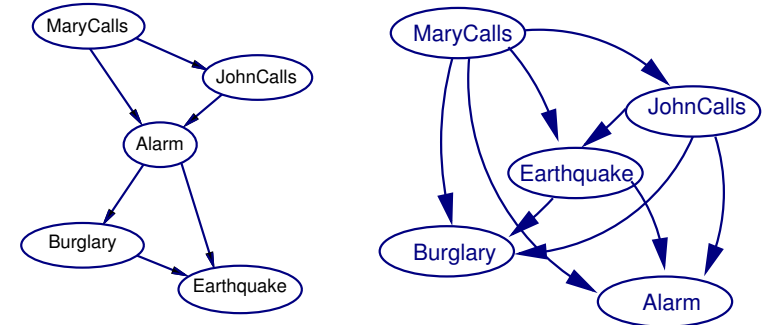
Belief Network Construction (cont'd)

The general procedure for Belief Network construction is as follows:

1. Choose the set of relevant variables X_i that describe the domain.
2. Choose an ordering of the variables.
3. While there are variables left:
 - (a) Pick a variable X_i and add a node to the network for that variable.
 - (b) Set $Parents(X_i)$ to some minimal set of nodes already in the net such that the conditional independence property is satisfied.
 - (c) Define the conditional probability table for node X_i .

58

Importance of Node Ordering in Belief Network Construction

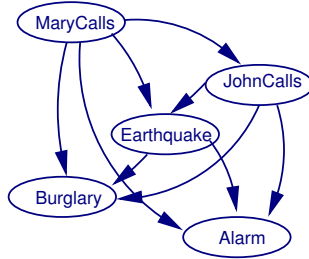


The resulting Belief Network can be vastly different when the order of insertion of nodes into the network is different.

- *MaryCalls, JohnCalls, Alarm, Burglary, Earthquake*
- *MaryCalls, JohnCalls, Earthquake, Burglary, Alarm*

60

Node Ordering and Joint Probability Tables



- Even with different graphs resulting from different node ordering, you can represent the same joint probability distribution.
- However, some represent the conditional independence relation much better than others.
- For example, the graph above requires the same number of entries as a full joint probability table.

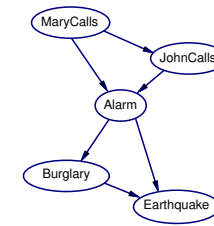
61

Strategy for Better Node Ordering

- For a node A to be a parent of another node B , A must be added to the network before B is added.
- Thus, a node that has a direct influence on other nodes should be added to the network first, i.e. add the **root causes** first.

63

Improper Node Ordering Can Cause Problems



- When adding nodes in the order of *MaryCalls*, *JohnCalls*, ..., conditional independence does not hold:

$$P(\text{JohnCalls}|\text{MaryCalls}) \neq P(\text{JohnCalls})$$
- Thus, *MaryCalls* has to become a parent of *JohnCalls*.
- This is because if Mary calls, it probably means that the alarm has gone off, so it makes it more likely that John calls.

62

Strategy for Better Node Ordering (cont'd)

When building the network, stick to a **causal model**, rather than the other way around (e.g. inferring cause given the effect).

Causal model is beneficial because of these reasons:

- conditional probability tables can be made smaller
- the conditional probabilities can be easier to come up with
- easier to reason about the domain using the network

64

Probabilistic Inference

- **Diagnostic inferences:** $P(\text{Cause}|\text{Effect})$
- **Causal inferences:** $P(\text{Effect}|\text{Cause})$
- **Intercausal inferences:** causes of a common effect (**explaining away**: cause has already been found)

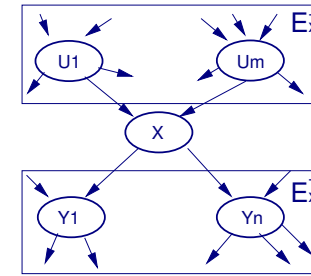
$$P(\text{Cause}|\text{Effect}) \gg P(\text{Cause}|\text{Effect} \wedge \text{OtherCause})$$

- **Mixed inferences:** combining two or more of the above

$$P(A|\text{CauseOfA} \wedge \text{EffectOfA})$$

65

Answering Queries: A Brief Outline



Given a set of evidence E , find the conditional probability $\mathbf{P}(X|E)$ where X is the query.

- Recursively determine the **causal support** E_X^+ for X .
- Recursively determine the **evidential support** E_X^- for X .

Note: This is only when the graph is **singly connected**.

66

Using Belief Networks and Probabilistic Inference

- Making decisions based on the derived probabilities and an agent's utility function.
- Deciding which additional variables to include in the model.
- Performing **sensitivity analysis** to find out which node is most important and thus should be more accurate.
- Explaining the results of reasoning.

67

Knowledge Engineering for Uncertain Reasoning

- Decide what to talk about (i.e. what to be included in the model). Gradually add more factors that can influence the current collection of events.
- Determine the variables to use and the range of the values.
- Encode general knowledge about the dependence between variables:
 1. qualitative: which variable depends on some other variable
 2. quantitative: probability value of the dependence (from experience, or from data gathered from a sample space)
- Encode a description of the specific problem instance: assign values to the variables.
- Pose queries to the inference procedure and get answers: what is the probability of X ? how sensitive are the values in the conditional probability tables to perturbations?

68

Key Points

- Constructing a belief network: what is the procedure? why does node ordering matter? how to order the nodes?
- Inference in belief networks: what are the kinds of inference? what is the general method?
- Knowledge engineering: how to formulate the idea and design a system.