# Overview

# 1.1 Acting under Uncertainty

## Logical agents versus decision-theoretic agents

Let $A$ be a proposition such as *heavy(Object1).*

For a logical agent $A$ is either false or true or  i.e.  either $A$ or $\mathrm{not}(A)$ is true.

When an agent knows enough facts about its environment, the logical approach enables it to derive plans that are guaranteed to work.

Unfortunately, agents almost never have access to the whole truth about their environment.

For a logical agent, it might be impossible to construct a complete and correct description of how its actions will work.

Agents must, therefore, live and act under uncertainty.

# Example

Agent wants to drive to the airport to catch a flight and is considering a plan, $A_{90}$, that involves leaving home 90 minutes before the flight departs and driving at a reasonable speed.

Even though the airport is only about 15 miles away, the agent will not be able to conclude with certainty that "Plan $A_{90}$ will get us to the airport in time."

# Example

Agent wants to drive to the airport to catch a flight and is considering a plan, $A_{90}$, that involves leaving home 90 minutes before the flight departs and driving at a reasonable speed.

Even though the airport is only about 15 miles away, the agent will not be able to conclude with certainty that "Plan $A_{90}$ will get us to the airport in time."

Instead, it reaches the weaker conclusion "Plan $A_{90}$ will get us to the airport in time,
- as long as the car doesn't break down
- or run out of gas,
- and we don't get into an accident,
- and there are no accidents on the bridge,
- and the plane doesn't leave early,
- and. . . ."

None of these conditions can be deduced, so the plan's success cannot be inferred.

# Example

Agent wants to drive to the airport to catch a flight and is considering a plan, $A_{90}$, that involves leaving home 90 minutes before the flight departs and driving at a reasonable speed.

Even though the airport is only about 15 miles away, the agent will not be able to conclude with certainty that "Plan $A_{90}$ will get us to the airport in time."

Instead, it reaches the weaker conclusion "Plan $A_{90}$ will get us to the airport in time,
- as long as the car doesn't break down
- or run out of gas,
- and we don't get into an accident,
- and there are no accidents on the bridge,
- and the plane doesn't leave early,
- and. . . ."

None of these conditions can be deduced, so the plan's success cannot be inferred.

This is an example of the qualification problem [1].

[1] The problem of reasoning about the conditions required for an event to have a given consequence is called the *qualification problem* and is of considerable interest to researchers working in default reasoning and nonmonotonic logic.

# Example (cont.)

The information, which the agent has, cannot guarantee any of these outcomes for $A_{90}$, but it can provide some degree of belief that they will be achieved.

Other plans, such as $A_{120}$, might increase the agent's belief that it will get to the airport on time, but also increase the likelihood of a long wait.

The right thing to do - *the rational decision* - therefore depends on both the relative importance of various goals and the likelihood that, and degree to which, they will be achieved.

# Handling uncertain knowledge

Rules for dental diagnosis using first-order logic

Consider the following rule:

$$\forall p \; Symptom \; (p, \; Toothache) \; \Rightarrow \; Disease \; (p, \; Cavity)$$

This rule is wrong. Not all patients with toothaches have cavities.

# Handling uncertain knowledge

Rules for dental diagnosis using first-order logic

Consider the following rule:

$$\forall p \; Symptom \; (p, \; Toothache) \; \Rightarrow \; Disease \; (p, \; Cavity)$$

This rule is wrong. Not all patients with toothaches have cavities.

$$\forall p \; Symptom \; (p, \; Toothache) \; \Rightarrow$$
$$Disease \; (p, \; Cavity) \lor Disease \; (p, \; GumDisease) \lor Disease \; (p, \; Abscess)...$$

To make the rule true, we have to add an almost unlimited list of possible causes.

# Handling uncertain knowledge

Rules for dental diagnosis using first-order logic

Consider the following rule:

$$\forall p \; Symptom \; (p, \; Toothache) \; \Rightarrow \; Disease \; (p, \; Cavity)$$

This rule is wrong. Not all patients with toothaches have cavities.

$$\forall p \; Symptom \; (p, \; Toothache) \; \Rightarrow$$
$$Disease \; (p, \; Cavity) \lor Disease \; (p, \; GumDisease) \lor Disease \; (p, \; Abscess)...$$

To make the rule true, we have to add an almost unlimited list of possible causes.

We could try turning the rule into a causal rule:

$$\forall p \; Disease(p, \; Cavity) \; \Rightarrow \; Symptom(p, \; Toothache)$$

this rule is not right either; not all cavities cause pain.

# Handling uncertain knowledge

Rules for dental diagnosis using first-order logic

Consider the following rule:

$$\forall p \; Symptom \; (p, \; Toothache) \; \Rightarrow \; Disease \; (p, \; Cavity)$$

This rule is wrong. Not all patients with toothaches have cavities.

$$\forall p \; Symptom \; (p, \; Toothache) \; \Rightarrow$$
$$Disease \; (p, \; Cavity) \vee Disease \; (p, \; GumDisease) \vee Disease \; (p, \; Abscess)...$$

To make the rule true, we have to add an almost unlimited list of possible causes.

We could try turning the rule into a causal rule:

$$\forall p \; Disease(p, \; Cavity) \; \Rightarrow \; Symptom(p, \; Toothache)$$

this rule is not right either; not all cavities cause pain.

The logical approach breaks down!

# Handling uncertain knowledge

Rules for dental diagnosis using first-order logic

Consider the following rule:

$$\forall p \; Symptom \; (p, \; Toothache) \; \Rightarrow \; Disease \; (p, \; Cavity)$$

This rule is wrong. Not all patients with toothaches have cavities.

$$\forall p \; Symptom \; (p, \; Toothache) \; \Rightarrow$$
$$Disease \; (p, \; Cavity) \lor Disease \; (p, \; GumDisease) \lor Disease \; (p, \; Abscess)...$$

To make the rule true, we have to add an almost unlimited list of possible causes.

We could try turning the rule into a causal rule:

$$\forall p \; Disease(p, \; Cavity) \; \Rightarrow \; Symptom(p, \; Toothache)$$

this rule is not right either; not all cavities cause pain.

The only way to fix the rule is to make it logically exhaustive: to augment the left-hand side with all the qualifications required for a cavity to cause a toothache.

IS THAT POSSIBLE???

*The logical approach breaks down!*

# Handling uncertain knowledge

Trying to use first-order logic to cope with a domain like medical diagnosis thus
fails for three main reasons:

*Laziness:*
It is too much work to list the complete set of antecedents or consequents needed
to ensure an exceptionless rule and too hard to use such rules.

*Theoretical ignorance:*
Medical science has no complete theory for the domain.

*Practical ignorance:*
Even if we know all the rules, we might be uncertain about a particular patient
because not all the necessary tests have been or can be run.

# Handling uncertain knowledge

Connection between toothaches and cavities is just not a logical consequence in either direction. Agent's knowledge can at best provide only a **degree of belief** in the relevant sentences.

Main tool for dealing with degrees of belief will be probability theory, which assigns to each sentence a numerical degree of belief between 0 and 1.

Assigning probability of 0 ~ unequivocal belief that the sentence is false.

Assigning probability of 1 ~ unequivocal belief that the sentence is true.

Probabilities between 0 and 1 correspond to intermediate degrees of belief in the truth of the sentence. The sentence itself is in fact either true or false.

Degree of belief is different from a degree of truth. A probability of 0.8 does not mean "80% true" but rather an 80% degree of belief (a fairly strong expectation).

# Handling uncertain knowledge

In probability theory, a sentence such as "The probability that the patient has a cavity is 0.8" is about the agent's beliefs, not directly about the world.

These beliefs depend on the percepts that the agent has received to date. As the agent receives new percepts, its probability assessments are updated to reflect the new evidence.

Before the evidence is obtained, we talk about prior or unconditional probability; after the evidence is obtained, we talk about posterior or conditional probability.

In most cases, an agent will have some evidence from its percepts and will be interested in computing the posterior probabilities of the outcomes it cares about.

# Uncertainty and rational decisions

Presence of uncertainty radically changes the way an agent makes decisions

Logical agent typically has a goal and executes any plan that is guaranteed to achieve it. An action can be selected or rejected on the basis of whether it achieves the goal, regardless of what other actions might achieve.

When uncertainty enters the picture, this is no longer the case.

Example: *$A_{90}$ plan for getting to the airport.*

Suppose it has a 95% chance of succeeding. Does this mean it is a rational choice?

Not necessarily: There might be other plans, such as $A_{120}$, with higher probabilities of success. If it is vital not to miss the flight, then it is worth risking the longer wait at the airport.

What about $A_{1440}$, a plan that involves leaving home 24 hours in advance?

To make such choices, an agent must first have preferences between the different possible outcomes of the various plans. A particular outcome is a completely specified state, including such factors as whether the agent arrives on time and the length of the wait at the airport.

# Utility theory

We will be using utility theory to represent and reason with preferences. Utility theory says that every state has a degree of usefulness, or utility, to an agent and that the agent will prefer states with higher utility [1] .

Preferences, as expressed by utilities, are combined with probabilities in the general theory of rational decisions called decision theory:

## Decision theory = probability theory + utility theory

The fundamental idea of decision theory is that an agent is rational if and only if it chooses the action that yields the highest expected utility, averaged over the possible outcomes of the action.

This is called the principle of Maximum Expected Utility (MEU).

---

[1] A utility function can even account for altruistic behavior, simply by including the welfare of others as one of the factors contributing to the agent's own utility.

# WHERE DO PROBABILITIES COME FROM?

## The Frequentist

The frequentist position is that the numbers can come only from experiments: if we test 100 people and find that 10 of them have a cavity, then we can say that the probability of a cavity is approximately 0.1. In this view, the assertion "the probability of a cavity is 0.1" means that 0.1 is the fraction that would be observed in the limit of infinitely many samples.

## The Objectivist

The objectivist view is that probabilities are real aspects of the universe - propensities of objects to behave in certain ways - rather than being just descriptions of an observer's degree of belief. For example, that a fair coin comes up heads with probability 0.5 is a propensity of the coin itself. In this view, the frequentist's measurements are attempts to observe these propensities.

## The Subjectivist

The subjectivist view describes probabilities as a way of characterizing an agent's beliefs, rather than as having any external physical significance.

# WHERE DO PROBABILITIES COME FROM?

## The Frequentist

The frequentist position is that the numbers can come only from experiments: if we test 100 people and find that 10 of them have a cavity, then we can say that the probability of a cavity is approximately 0.1. In this view, the assertion "the probability of a cavity is 0.1" means that 0.1 is the fraction that would be observed in the limit of infinitely many samples.

7 5 1 4 8

## The Objectivist

The objectivist view is that probabilities are real aspects of the universe - propensities of objects to behave in certain ways - rather than being just descriptions of an observer's degree of belief. For example, that a fair coin comes up heads with probability 0.5 is a propensity of the coin itself. In this view, the frequentist's measurements are attempts to observe these propensities.

1€

## The Subjectivist

The subjectivist view describes probabilities as a way of characterizing an agent's beliefs, rather than as having any external physical significance.

In the end, even a strict frequentist position involves subjective analysis, so the difference probably has little practical importance.

# 1.2 Basic Probability Notation

## Introduction

‣ Probability theory uses an extension of propositional logic for its sentences

‣ Language is slightly more expressive than propositional logic

‣ Dependence on experience is reflected in the syntactic distinction between prior probability statements, which apply before any evidence is obtained, and conditional probability statements, which include the evidence explicitly

## Propositions

‣ Basic elements: **random variable**
(referring to a "part" of the world whose "status" is initially unknown)

◊ Boolean random variables: *Cavity* with domain ⟨*true, false*⟩
i.e. value(Cavity) ∈ ⟨*true, false*⟩; *Cavity = true ~ cavity, Cavity = false ~ ¬cavity*

◊ Discrete random variables: with countable domains, e.g. *Weather* with
⟨*sunny, rainy, cloudy, snow*⟩

◊ Continuous random variables: domain can be reals, e.g. $X = 4.02$ or $X \geq 4.02$
or subset of intervals e.g. [0,1]

Elementary propositions such as *Cavity = true* and *Toothache = false*, can be combined to form complex propositions using all the standard logical connectives: e.g. *Cavity = true ∧ Toothache = false* or *cavity ∧ ¬toothache*

# Atomic events

▸ An atomic event is complete specification of the state of the world about which the agent is uncertain {*Cavity, Toothache, …*}

▸ Assignment of particular values to all the variables of which the world is composed e.g. *Cavity = false ∧ Toothache = true*

▸ Atomic events have some important properties:
  • mutually exclusive: *cavity ∧ toothache* and *cavity ∧ ¬ toothache*
  • set of all possible atomic events is exhaustive - at least one must be the case
  • any proposition is logically equivalent to the disjunction of all atomic events that entail the truth of the proposition:
    *cavity* is equivalent to (*cavity ∧ toothache*) ∨ (*cavity ∧ ¬ toothache*)

# Prior probability

▸ Unconditional or prior probability associated with a proposition $a$ is the degree of belief accorded to it in the absence of any other information written as $P(a)$.

Example: if the prior probability that I have a cavity is 0.1

$$P(Cavity = true) = 0.1 \quad \text{or} \quad P(cavity) = 0.1$$

It is important to remember that $P(a)$ can be used only when there is no other information. As soon as some new information is known, we must reason with the conditional probability of a given that new information.

▸ Notation

$$P(\ Weather = sunny\ ) = 0.7$$
$$P(\ Weather = rain\ ) = 0.2$$
$$P(\ Weather = cloudy\ ) = 0.08$$
$$P(\ Weather = snow\ ) = 0.02$$

is often abbreviated as

$$\mathbf{P}(\ Weather\ ) = \langle 0.7, 0.2, 0.08, 0.02 \rangle$$

▸ $\mathbf{P}(\ Weather)$ describes the prior probability distribution for the random variable $Weather$

▸ **P**( *Weather, Cavity*) denotes probabilities of all combinations of the values of a set of random variables

▸ and is represented by a 4 x 2 table of probabilities called the joint probability distribution of *Weather* and *Cavity*

▸ Sometimes it will be useful to think about the complete set of random variables used to describe the world: a joint probability distribution that covers this complete set is called the **full joint probability distribution**

Example: if the world consists of just the variables *Cavity, Toothache,* and *Weather*, then the full joint distribution is given by

     **P**(*Cavity, Toothache, Weather*)

This joint distribution can be represented as a 2 x 2 x 4 table with 16 entries.

A full joint distribution specifies the probability of every atomic event and is therefore a complete specification of one's uncertainty about the world in question.

▸ For continuous variables, it is not possible to write out the entire distribution as a table, because there are infinitely many values.

▸ Instead, one usually defines the probability that a random variable takes on some value $x$ as a parameterized function of $x$.

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}}\exp(-\frac{(x-\mu)^2}{2\sigma^2})$$

▸ Probability distributions for continuous variables are called probability density functions.

Density functions differ in meaning from discrete distributions.

How???

# Conditional probability

▸ Once agent obtained some evidence concerning the previously unknown random variables making up the domain, prior probabilities are no longer applicable. Instead, we use conditional or posterior probabilities.

▸ "the probability of $a$, given that all we know is $b$." The notation used is

$$P(a|b),$$

where $a$ and $b$ are any propositions

▸ Example: $P(cavity \mid toothache) = 0.8$

indicates that if a patient is observed to have a toothache and no other information is yet available, then the probability of the patient's having a cavity will be 0.8.

▸ Conditional probabilities can be defined in terms of unconditional probabilities

$$P(a \mid b) = \frac{P(a \wedge b)}{P(b)} \qquad\qquad (1.1)$$

which holds whenever $P(b) > 0$. This equation can also be written as

$$P(a \wedge b) = P(a \mid b)P(b)$$

which is called the product rule.

▸ We can also have it the other way around:

$$P(a \wedge b) = P(b \mid a)P(a)$$

▸ We can use the **P notation for conditional distributions**. $\mathbf{P}(X \mid Y)$ gives the values of $P(X = x_i \mid Y = y_i)$ for each possible $i, j$.

$$P(X = x_1 \wedge Y = y_1) = P(X = x_1 \mid Y = y_1)P(Y = y_1)$$
$$P(X = x_1 \wedge Y = y_2) = P(X = x_1 \mid Y = y_2)P(Y = y_2)$$
$$\vdots$$

We can combine all these into the single equation

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

This denotes a set of equations relating the corresponding individual entries in the tables, not a matrix multiplication of the tables.

---

The "|" operator has the lowest possible precedence, so $P(a \wedge b \mid c \vee d)$ means $P((a \wedge b) \mid (c \vee d))$

▸ It is tempting, but wrong, to view conditional probabilities as if they were logical implications with uncertainty added. For example, the sentence $P(a \mid b) = 0.8$ cannot be interpreted to mean "whenever $b$ holds, conclude that $P(a)$ is 0.8."

▸ Such an interpretation would be wrong on two counts:

  Why?

▸ It is tempting, but wrong, to view conditional probabilities as if they were logical implications with uncertainty added. For example, the sentence $P(a \mid b) = 0.8$ cannot be interpreted to mean "whenever $b$ holds, conclude that $P(a)$ is 0.8."

▸ Such an interpretation would be wrong on two counts:

- $P(a)$ always denotes the prior probability of a, not the posterior probability given some evidence;

- statement $P(a \mid b) = 0.8$ is relevant just when $b$ is the only available evidence. When additional information $c$ is available, the degree of belief in $a$ is $P(a \mid b \wedge c)$ which may have little relation to $P(a \mid b)$

# 1.3    The Axioms of Probability

▸   We have defined a syntax for propositions and for prior and conditional probability statements about those propositions

▸   Now we must provide some sort of semantics for probability statements.

Basic axioms defining probability scale and end points

1.  All probabilities are between 0 and 1. For any proposition $a$,

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

2.  Necessarily true (i.e., valid) propositions have probability 1, and necessarily false (i.e., unsatisfyable) propositions have probability 0.

$$P(true) = 1 \qquad P(false) = 0$$

Axiom connecting the probabilities of logically related propositions

3. The probability of a disjunction is given by

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

These three axioms are often called Kolmogorov's axioms in honor of the Russian mathematician Andrei Kolmogorov, who showed how to build up the rest of probability theory from this simple foundation.

# Using the axioms of probability

‣ We can derive a variety of useful facts from the basic axioms. For example, the familiar rule for negation follows by substituting $\neg a$ for $b$ in axiom 3, giving us:

$$P(a \vee \neg a) =$$

# Using the axioms of probability

▸ We can derive a variety of useful facts from the basic axioms. For example, the familiar rule for negation follows by substituting $\neg a$ for $b$ in axiom 3, giving us:

$$P(a \vee \neg a) = P(a) + P(\neg a) - P(a \wedge \neg a) \quad \text{(by axiom 3 with } b = \neg a\text{)}$$

$$P(true) = P(a) + P(\neg a) - P(false) \quad \text{(by logical equivalence)}$$

$$1 = P(a) + P(\neg a) \quad \text{(by axiom 2)}$$

$$P(\neg a) = 1 - P(a) \quad \text{(by algebra).}$$

▸ The third line of this derivation is itself a useful fact and can be extended from the Boolean case to the general discrete case.

Let the discrete variable $D$ have the domain $\langle d_1, ..., d_n \rangle$. Then it is easy to show that

$$\sum_{i=1}^{n} P(D = d_i) = 1$$

That is, any probability distribution on a single variable must sum to 1. Also, any joint probability distribution on any set of variables must sum to 1.

For continuous variables, the summation is replaced by an integral

$$\int_{-\infty}^{\infty} P(X = x)\, dx = 1$$

▸ From axiom 3 we can further derive the following relation

The probability of a proposition is equal to the sum of the probabilities of the atomic events in which it holds

$$P(a) = \sum_{e_i \in e(a)} P(e_i) \qquad\qquad (1.2)$$

Equation provides a simple method for computing the probability of any proposition given a full joint distribution that specifies the probabilities of all atomic events

# Why the axioms of probability are reasonable

▸ Axioms of probability can be seen as restricting the set of probabilistic beliefs that an agent can hold. This is somewhat analogous to the logical case, where a logical agent cannot simultaneously believe $A$, $B$, and $\neg (A \wedge B)$, for example.

▸ With probabilities, statements refer not to the world directly, but to the agent's own state of knowledge. Why, then, can an agent not hold the following set of beliefs, which clearly violates axiom 3?

$$
\begin{array}{lll}
P(a) = 0.4 & P(a \wedge b) = 0.0 & \qquad (1.3) \\
P(b) = 0.3 & P(a \vee b) = 0.8 &
\end{array}
$$

# Why the axioms of probability are reasonable

▸ Axioms of probability can be seen as restricting the set of probabilistic beliefs that an agent can hold. This is somewhat analogous to the logical case, where a logical agent cannot simultaneously believe $A$, $B$, and $\neg (A \wedge B)$, for example.

▸ With probabilities, statements refer not to the world directly, but to the agent's own state of knowledge. Why, then, can an agent not hold the following set of beliefs, which clearly violates axiom 3?

$$P(a) = 0.4 \qquad \mathrm{P}(a \wedge b) = 0.0 \qquad\qquad (1.3)$$
$$P(b) = 0.3 \qquad \mathrm{P}(a \vee b) = 0.8$$

▸ This question has been the subject of decades of intense debate between those who advocate the use of probabilities as the only legitimate form for degrees of belief and those who advocate alternative approaches.

▸ Here, we give one argument for the axioms of probability, first stated in 1931 by Bruno de Finetti.

# Finetti's Argument

▸ The key to de Finetti's argument is the connection between degree of belief and actions.

▸ If an agent has some degree of belief in a proposition $a$, then the agent should be able to state odds at which it is indifferent to bet for or against $a$.

Think of it as a game between two agents: Agent 1 states "my degree of belief in event a is 0.4." Agent 2 is then free to choose whether to bet for or against $a$, at stakes that are consistent with the stated degree of belief.

Agent 2 can choose to
     bet $4 against Agent 1's $6 that a will occur or
     bet $6 against Agent 1's $4 that a will not occur.

i.e. Agent 1 bets $4 that a will occur,
Agent 2 bets $6 that a will not occur

If Agent 1 expresses a set of degrees of belief that violate the axioms of probability theory then there is a combination of bets of Agent 2 that guarantees that Agent 1 will lose money every time.

## Finetti's Argument (continued)

Suppose that Agent 1 has the set of inconsistent beliefs (see Equation (1.3) and Figure)

If Agent 2 chooses to bet $4 on $a$ , $3 on $b$, and $2 on $\neg(a \vee b)$ then Agent 1 always loses money, regardless of the outcomes for $a$ and $b$.

| Agent1 | | Agent2 | | Outcome for Agent 1 | | | |
|--------|--------|--------|--------|--------|--------|--------|--------|
| Proposition | Belief | Bet | Stakes | $a \wedge b$ | $a \wedge \neg b$ | $\neg a \wedge b$ | $\neg a \wedge \neg b$ |
| $a$ | 0.4 | $a$ | 4 to 6 | -6 | -6 | 4 | 4 |
| $b$ | 0.3 | $b$ | 3 to 7 | -7 | 3 | -7 | 3 |
| $a \vee b$ | 0.8 | $\neg (a \vee b)$ | 2 to 8 | 2 | 2 | 2 | -8 |
| | | | | -11 | -1 | -1 | -1 |

# 1.4    Inference Using Full Joint Distributions

‣ We will describe a simple method for <mark>probabilistic inference</mark> (computation from observed evidence of posterior probabilities for query propositions)

‣ We will use the full joint distribution as the "knowledge base" from which answers to all questions may be derived

‣ Simple example: a domain consisting of just the three Boolean variables *Toothache, Cavity*, and *Catch*

query: the dentist's nasty steel probe catches in my tooth

Full joint distribution is a 2 x 2 x 2 table as shown in the following figure

| | *toothache* | | $\neg$ *toothache* | |
|---|---|---|---|---|
| | *catch* | $\neg$ *catch* | *catch* | $\neg$ *catch* |
| *cavity* | 0.108 | 0.012 | 0.072 | 0.008 |
| $\neg$ *cavity* | 0.016 | 0.064 | 0.144 | 0.576 |
| A full joint distribution for the *Toothache, Cavity, Catch* world. | | | | |

|  | *toothache* | | ¬ *toothache* | |
|  | *catch* | ¬ *catch* | *catch* | ¬ *catch* |
|---|---|---|---|---|
| *cavity* | 0.108 | 0.012 | 0.072 | 0.008 |
| ¬ *cavity* | 0.016 | 0.064 | 0.144 | 0.576 |
| A full joint distribution for the *Toothache, Cavity, Catch* world. | | | | |

- Notice that the probabilities in the joint distribution sum to 1, as required by the axioms of probability.

- Notice also that Equation (1.2) gives us a direct way to calculate the probability of any proposition, simple or complex: We simply identify those atomic events in which the proposition is true and add up their probabilities.

- For example, there are six atomic events in which *cavity* ∨ *toothache* holds:

  $P(cavity \lor tootache) =$ ???

| | *toothache* | | ¬ *toothache* | |
|---|---|---|---|---|
| | *catch* | ¬ *catch* | *catch* | ¬ *catch* |
| *cavity* | 0.108 | 0.012 | 0.072 | 0.008 |
| ¬ *cavity* | 0.016 | 0.064 | 0.144 | 0.576 |
| A full joint distribution for the *Toothache, Cavity, Catch* world. | | | | |

▸ Notice that the probabilities in the joint distribution sum to 1, as required by the axioms of probability.

▸ Notice also that Equation (1.2) gives us a direct way to calculate the probability of any proposition, simple or complex: We simply identify those atomic events in which the proposition is true and add up their probabilities.

▸ For example, there are six atomic events in which *cavity* ∨ *toothache* holds:

   *P(cavity* ∨ *tootache)* = 0.108 + 0.012 + 0.072 + 0.008 + 0.016 + 0.064 = 0.28

‣ One particularly common task is to extract the distribution over some subset of variables or a single variable.

‣ For example, adding the entries in the first row gives the <mark>unconditional or marginal probability</mark> of *cavity*:

$$P(cavity) = 0.108 + 0.012 + 0.072 + 0.008 = 0.2.$$

This process is called <mark>marginalization, or summing out</mark> - because the variables other than *Cavity* are summed out.

‣ We can write the following general marginalization rule for any sets of variables $Y$ and $Z$:

$$\mathbf{P(Y)} = \sum_z \mathbf{P(Y,z)} \qquad (1.4)$$

‣ A variant of this rule involves conditional probabilities instead of joint probabilities, using the product rule:

$$\mathbf{P(Y)} = \sum_z \mathbf{P(Y\,|\,z)}P(z) \qquad (1.5)$$

This rule is called conditioning.

Marginalization and conditioning will turn out to be useful rules for all kinds of derivations involving probability expressions.

▸ Often we will be interested in computing conditional probabilities of some variables, given evidence about others.

▸ Example: probability of cavity, given evidence of a toothache

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

$$= \frac{0.108 + 0.012}{0.108 + 0.012 + 0.016 + 0.064} = 0.6$$

probability that there is no cavity, given a toothache

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

$$= \frac{0.016 + 0.064}{0.108 + 0.012 + 0.016 + 0.064} = 0.4$$

▸ Notice that in these two calculations the term $1/P(toothache)$ remains constant, no matter which value of $Cavity$ we calculate; can be viewed as a normalization constant for the distribution $\mathbf{P}(Cavity|toothache)$, ensuring that it adds up to 1.

▸ We will use $\alpha$ to denote such constants and write the two preceding equations in one:

$\mathbf{P}(Cavity \mid toothache) = \alpha\,\mathbf{P}(Cavity,\ toothache)$

$\qquad = \alpha\,[\mathbf{P}(Cavity,\ toothache,\ catch) + \mathbf{P}(Cavity,\ toothache,\ \neg\ catch)]$

$\qquad = \alpha\,[\langle\,0.108, 0.016\,\rangle + \langle\,0.012, 0.064\,\rangle] = \alpha\,\langle\,0.12, 0.08\,\rangle = \langle\,0.6, 0.4\,\rangle\,.$

Normalization will turn out to be a useful shortcut in many probability calculations.

General inference procedure:

Let $X$ be the query variable (*Cavity* in the example), let $\mathbf{E}$ be the set of evidence variables (just *Toothache* in the example), let $\mathbf{e}$ be the observed values for them, and let $\mathbf{Y}$ be the remaining unobserved variables (just *Catch* in the example). The query is $\mathbf{P}(X \,|\, \mathbf{e})$ and can be evaluated as

$$\mathbf{P}(X \,|\, \mathbf{e}) = \alpha\, \mathbf{P}(X, \mathbf{e}) = \alpha \sum_{y} \mathbf{P}(X, \mathrm{e}, \mathrm{y})$$

<u>Notice</u> together the variables $\mathrm{X}$, $\mathbf{E}$, and $\mathbf{Y}$ constitute the complete set of variables for the domain, so $\mathbf{P}(\mathrm{X}, \mathrm{e}, \mathrm{y})$ is simply a subset of probabilities from the full joint distribution

# 1.5  Independence

▸ Let us expand the full joint distribution in Figure 13.3 by adding a fourth variable, *Weather*. The full joint distribution then becomes

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

which has 32 entries (because *Weather* has four values).

▸ What relationship do these additions have to each other and to the original three-variable table??

For example, how are $P(\textit{toothache, catch, cavity, Weather} = \textit{cloudy})$ and $P(\textit{toothache, catch, cavity})$ related?

to answer this question we use the product rule:

$$P(\textit{toothache, catch, cavity, Weather} = \textit{cloudy}) =$$
$$= P(\textit{Weather} = \textit{cloudy} \mid \textit{toothache, catch, cavity}) \, P(\textit{toothache, catch, cavity}) \,.$$

▸ One's dental problems typically do not influence the weather. Therefore,

$$P(\textit{Weather} = \textit{cloudy} \mid \textit{toothache, catch, cavity}) = P(\textit{Weather} = \textit{cloudy}) \quad (1.7)$$
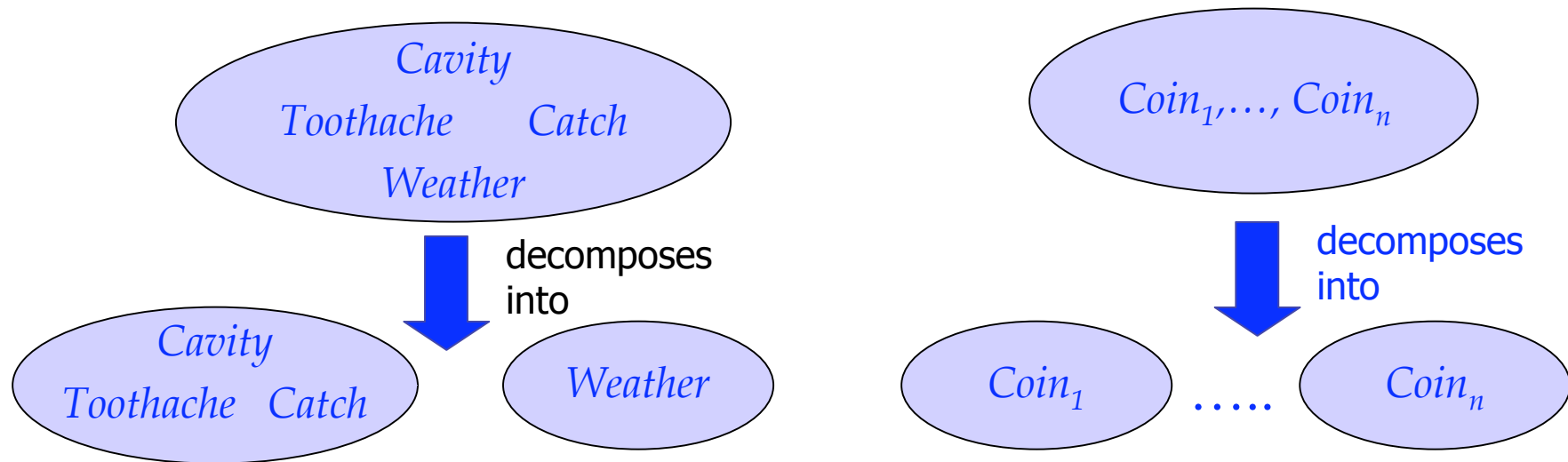
▸ From this, we can deduce

$$P(\textit{toothache, catch, cavity, Weather} = \textit{cloudy}) =$$
$$= P(\textit{Weather} = \textit{cloudy})\ P(\textit{toothache, catch, cavity})$$

▸ A similar equation exists for every entry in **P**(*Toothache, Catch, Cavity, Weather*). we can write the general equation

**P**(*Toothache, Catch, Cavity, Weather*) = **P**(*Toothache, Catch, Cavity*)**P**(*Weather*) .

Thus, the 32-element table for four variables can be constructed from one 8-element table and one four-element table



Two examples of factoring a large joint distribution into smaller distributions, using absolute independence (weather and dental problems are independent and so are coin flips)

▸ The property we used in writing Equation (1.7) is called independence (also marginal independence and absolute independence).

▸ Independence between propositions a and b can be written as

$$P(a|b) = P(a) \quad \text{or} \quad P(b|a) = P(b) \quad \text{or} \quad P(a \wedge b) = P(a)P(b) \qquad (1.8)$$

All these forms are equivalent.

▸ Independence between variables $X$ and $Y$ can be written as follows:

$$\mathbf{P}(X|Y) = \mathbf{P}(X) \quad \text{or} \quad \mathbf{P}(Y|X) = \mathbf{P}(Y) \quad \text{or} \quad \mathbf{P}(X \wedge Y) = \mathbf{P}(X)\mathbf{P}(Y)$$

Independence assertions are usually based on knowledge of the domain.

They can dramatically reduce the amount of information necessary to specify the full joint distribution.

If the complete set of variables can be divided into independent subsets, then the full joint can be factored into separate joint distributions on those subsets.

# 1.6  Bayes' Rule and Its Use

▸ Product rule can be written in two forms because of the commutativity of conjunction:

$$P(a \wedge b) = P(a \mid b)P(b)$$

$$P(a \wedge b) = P(b \mid a)P(a)$$

▸ Equating the two right-hand sides and dividing by $P(a)$, we get

$$P(b \mid a) = \frac{P(a \mid b)P(b)}{P(a)}$$

known as Bayes' rule (also Bayes' law or Bayes' theorem).

▸ This simple equation underlies all modern AI systems for probabilistic inference.

▸ The more general case of multi-valued variables can be written in the P notation (probability distribution) as

$$\mathbf{P}(Y \mid X) = \frac{\mathbf{P}(X \mid Y)\mathbf{P}(Y)}{\mathbf{P}(X)} \quad \text{or more general} \quad \mathbf{P}(Y \mid X, \mathbf{e}) = \frac{\mathbf{P}(X \mid Y, \mathbf{e})\mathbf{P}(Y \mid \mathbf{e})}{\mathbf{P}(X \mid \mathbf{e})}$$

conditionalized on some background evidence $\mathbf{e}$; where again this is to be taken as representing a set of equations, each dealing with specific values of the variables.

# Applying Bayes' rule: The simple case

▸ On the surface, Bayes' rule does not seem very useful. It requires three terms - a conditional probability and two unconditional probabilities - just to compute one conditional probability

▸ Bayes' rule is useful in practice, because there are many cases where we do have good probability estimates for these three numbers and need to compute the fourth.

## Example

In medical diagnosis, we often have conditional probabilities on causal relationships and want to derive a diagnosis.

- A doctor knows that the disease meningitis causes the patient to have a stiff neck, say, 50% of the time.

- The doctor also knows some unconditional facts: the prior probability that a patient has meningitis is 1/50000, and the prior probability that any patient has a stiff neck is 1/20.

- ...

## Example (cont.)

- Letting $s$ be the proposition that the patient has a stiff neck and $m$ be the proposition that the patient has meningitis, we have

$$P(s \,/\, m) = 0.5$$

$$P(m) = 1 \,/\, 50000$$

$$P(s) = 1 \,/\, 20$$

$$P(m \,/\, s) = \frac{P(s \mid m)P(m)}{P(s)} = \frac{0.5 \times 1 \,/\, 50000}{1 \,/\, 20} = 0.0002$$

That is we expect only 1 in 5000 patients with a stiff neck to have meningitis.

diagnostic inference:

$P(m \mid s) \sim$ probability of meningitis having evidence on a stiff neck

Question: Why would one have the conditional probability available in one direction, but not the other?

- perhaps doctor knows that stiff neck implies meningitis in 1 out of 5000 cases;

- that is, the doctor has quantitative information in the diagnostic direction i.e.from symptoms to causes  → no need to use Bayes' rule.

Unfortunately, diagnostic knowledge can be more fragile than causal knowledge.

- if there is a sudden epidemic of meningitis, the unconditional probability of meningitis, $P(m)$, will go up.

- doctor deriving diagnostic probability $P(m \mid s)$ directly from statistical observation of patients before the epidemic will have no idea how to update the value

- but doctor who computes $P(m \mid s)$ from the other three values will see that $P(m \mid s)$ should go up proportionately with $P(m)$.

- Most importantly, causal information $P(s \mid m)$ is *unaffected* by the epidemic

The use of this kind of direct causal or model-based knowledge provides the crucial robustness needed to make probabilistic systems feasible in the real world.

## Using Bayes' rule: Combining evidence

‣ Probabilistic information is often available in the form $P(\mathit{effect} \mid \mathit{cause})$.

‣ What, if we have two or more pieces of evidence, e.g.
$$P(\mathit{Cavity} \mid \mathit{toothache} \wedge \mathit{catch}) \text{ ?}$$

‣ Given, the full joint distribution one can read off the answer:

$\mathbf{P}(\mathit{Cavity} \mid \mathit{toothache} \wedge \mathit{catch}) = \alpha \langle\, 0.108, 0.016 \,\rangle \approx \langle\, 0.871, 0.129 \,\rangle$

however, that such an approach will not scale up to larger numbers of variables.

‣ We can try using use Bayes' rule to reformulate the problem
$\mathbf{P}(\mathit{Cavity} \mid \mathit{toothache} \wedge \mathit{catch}) = \alpha\, \mathbf{P}(\mathit{toothache} \wedge \mathit{catch} \mid \mathit{Cavity})\, \mathbf{P}(\mathit{Cavity})$ (1.12)

For this we need to know the conditional probabilities of the conjunction $\mathit{toothache} \wedge \mathit{catch}$ for each value of $\mathit{Cavity}$. That might be feasible for just two evidence variables, but again it will not scale up
($n$ possible evidence variables $2^{n}$ possible combinations of observed values)
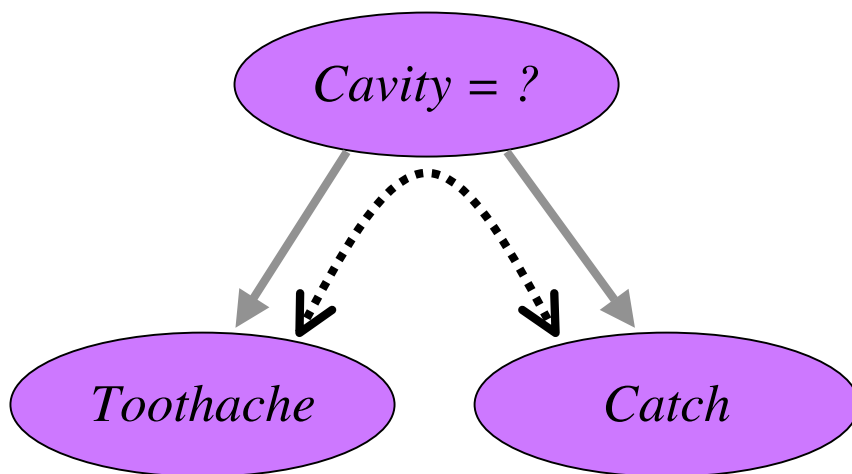
So what can we do???

# Conditional independence

▸ It would be nice if *Toothache* and *Catch* were independent, but they are not: if the probe catches in the tooth, it probably has a cavity and that probably causes a toothache.
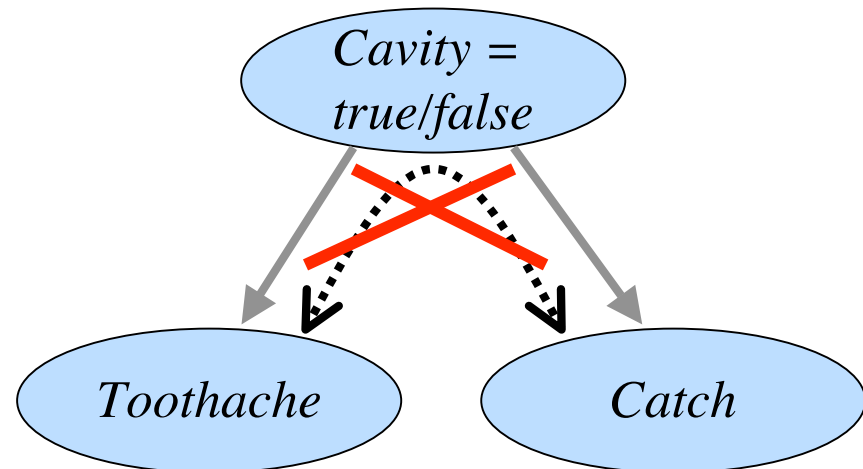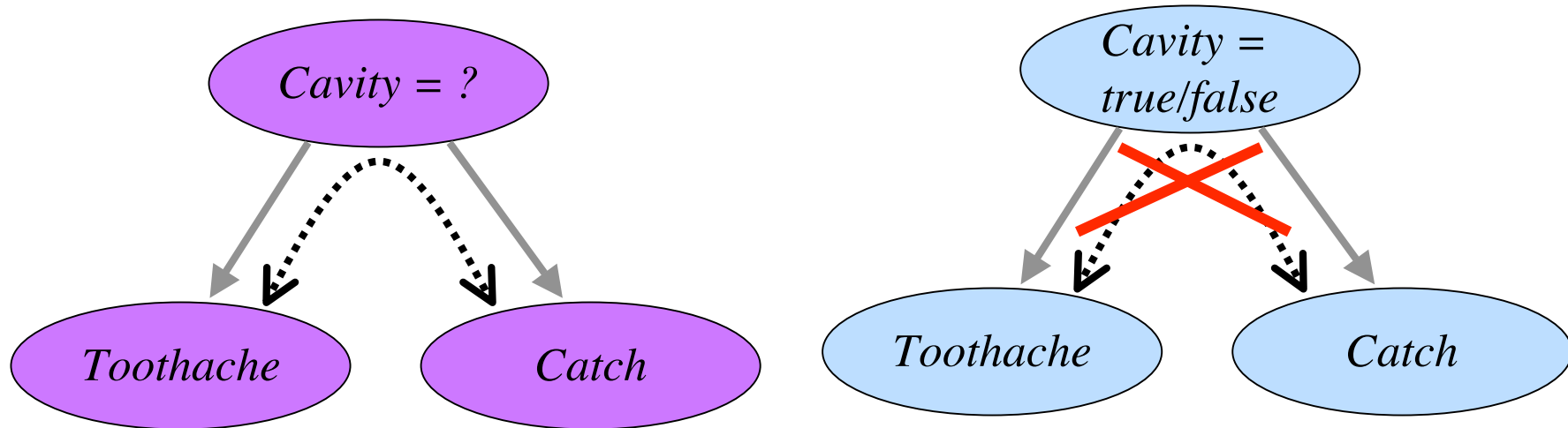


Without evidence on *Cavity* our beliefs regarding *Toothache* and *Catch* may affect each other

# Conditional independence

▸ It would be nice if *Toothache* and *Catch* were independent, but they are not: if the probe catches in the tooth, it probably has a cavity and that probably causes a toothache.

▸ Variables are independent, however, given presence or absence of a cavity. Each is directly caused by the cavity, but neither has a direct effect on the other.



Without evidence on *Cavity* our beliefs regarding *Toothache* and *Catch* may affect each other

Given evidence on *Cavity* (our beliefs regarding) *Toothache* and *Catch* have no affect each other !!!

▸ Mathematically, this property is written as

$\mathbf{P}(\textit{toothache} \wedge \textit{catch} \mid \textit{Cavity}) = \mathbf{P}(\textit{toothache}\mid\textit{Cavity}) \, \mathbf{P}(\textit{catch}\mid\textit{Cavity})$   (1.13)

conditional independence of *toothache* and *catch* given *Cavity*.

▸ We can plug it into Equation (1.12) to obtain the probability of a cavity:

$\mathbf{P}(\textit{Cavity} \mid \textit{toothache} \wedge \textit{catch}) = \alpha \, \mathbf{P}(\textit{toothache} \mid \textit{Cavity}) \, \mathbf{P}(\textit{catch} \mid \textit{Cavity}) \, \mathbf{P}(\textit{Cavity})$ (1.14)

▸ The general definition of conditional independence of two variables $X$ and $Y$, given a third variable $Z$ is

$\mathbf{P}(X, Y \mid Z) = \mathbf{P}(X \mid Z) \, \mathbf{P}(Y \mid Z)$

▸ Absolute independence assertions allow a decomposition of the full joint distribution into much smaller pieces.

▸ The same is true for conditional independence assertions: given the assertion in (1.14), we can derive a decomposition as follows:

$\mathbf{P}(Cavity, Catch, Toothache)$

    $= \mathbf{P}(Toothache, Catch \mid Cavity)\,\mathbf{P}(Cavity)$     (product rule)

    $= \mathbf{P}(Toothache \mid Cavity)\,\mathbf{P}(Catch \mid Cavity)\,\mathbf{P}(Cavity)$   (using 1.14)

▸ Original large table is decomposed into three smaller tables:

original table has seven independent numbers ($2^3$ - 1). The smaller tables contain five independent numbers ($2 \times (2^1$ - 1) for each conditional probability distribution and $2^1$ - 1 for the prior on *Cavity*).

▸ might not seem to be major triumph, but for $n$ symptoms that are all conditionally independent given *Cavity*, the size of the representation grows as $O(n)$ instead of $O(2^n)$

▸ Conceptually, *Cavity* separates *Toothache* and *Catch* because it is a direct cause of both of them (*d-separation*)
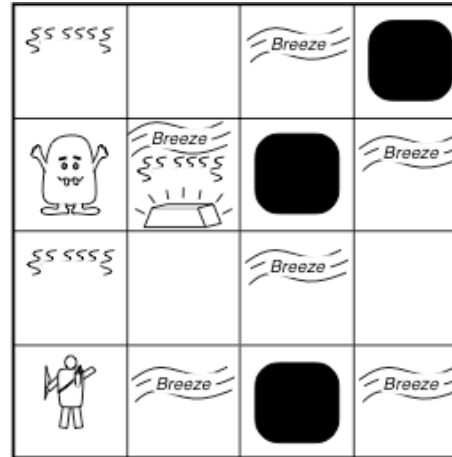
# Conditional independence

- Conditional independence assertions can allow probabilistic systems to scale up

- They are much more commonly available than absolute independence assertions.

- Decomposition of large probabilistic domains into weakly connected subsets via conditional independence is one of the most important developments in the recent history of AI.

▸ Commonly occurring pattern in which a single cause directly influences a number of effects, all of which are conditionally independent, given the cause. The full joint distribution can be written as

$$\mathbf{P}(Cause, Effect_1 \ldots, Effect_n) = \mathbf{P}(Cause) \prod_i \mathbf{P}(Effect_i \mid Cause) \, .$$

# 1.7 The Wumpus World

▸ The wumpus world is a cave consisting of rooms connected by passageways.



▸ Lurking somewhere in the cave is the wumpus, a beast that eats anyone who enters its room.

▸ Wumpus can be shot by an agent, but the agent has only one arrow.

▸ Some rooms contain bottomless pits that will trap anyone who wanders into these rooms (except for the wumpus).

▸ The only mitigating feature of living in this environment is the possibility of finding a heap of gold.

Wumpus world makes an excellent test bed environment for intelligent agents.

◇ **Performance measure:** + 1000 for picking up the gold, -1000 for falling into a pit or being eaten by the wumpus, -1 for each action taken and -10 for using up the arrow.

◇ **Environment:** 4 x 4 grid of rooms; agent always starts in the square labeled [1,1], facing to the right; locations of the gold and the wumpus are chosen randomly, with a uniform distribution, from the squares other than the start square; in addition, each square other than the start can be a pit, with probability 0.2.

◇ **Actuators:** The agent can *move forward*, *turn left by* 90°, or *turn right by* 90°. The agent dies if it enters a square containing a pit or a live wumpus.

Moving forward has no effect if there is a wall in front of the agent.

The action *grab* can be used to pick up an object that is in the same square as the agent.

The action *shoot* can be used to fire an arrow in a straight line in the direction the agent is facing. Arrow either hits (and hence kills) wumpus or hits wall. The agent only has one arrow, so only the first shoot action has any effect.

◇ **Sensors:** The agent has five sensors, each of which gives a single bit of information:
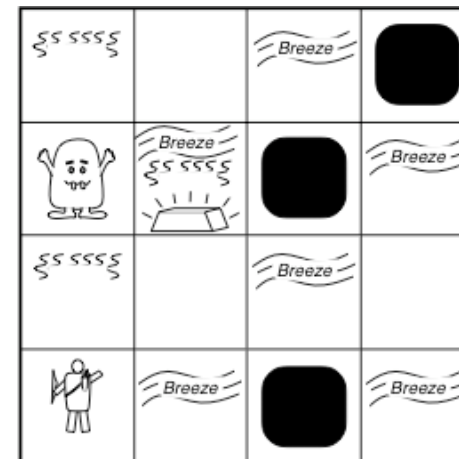
- in the square containing the wumpus and in the directly (not diagonally) adjacent squares the agent will perceive a stench,
- in the squares directly adjacent to a pit, the agent will perceive a breeze,
- in the square where the gold is, the agent will perceive a glitter,
- when an agent walks into a wall, it will perceive a bump,
- when the wumpus is killed, it emits a woeful scream that can be perceived anywhere in the cave.

Percepts will be given to agent as a list of five symbols:

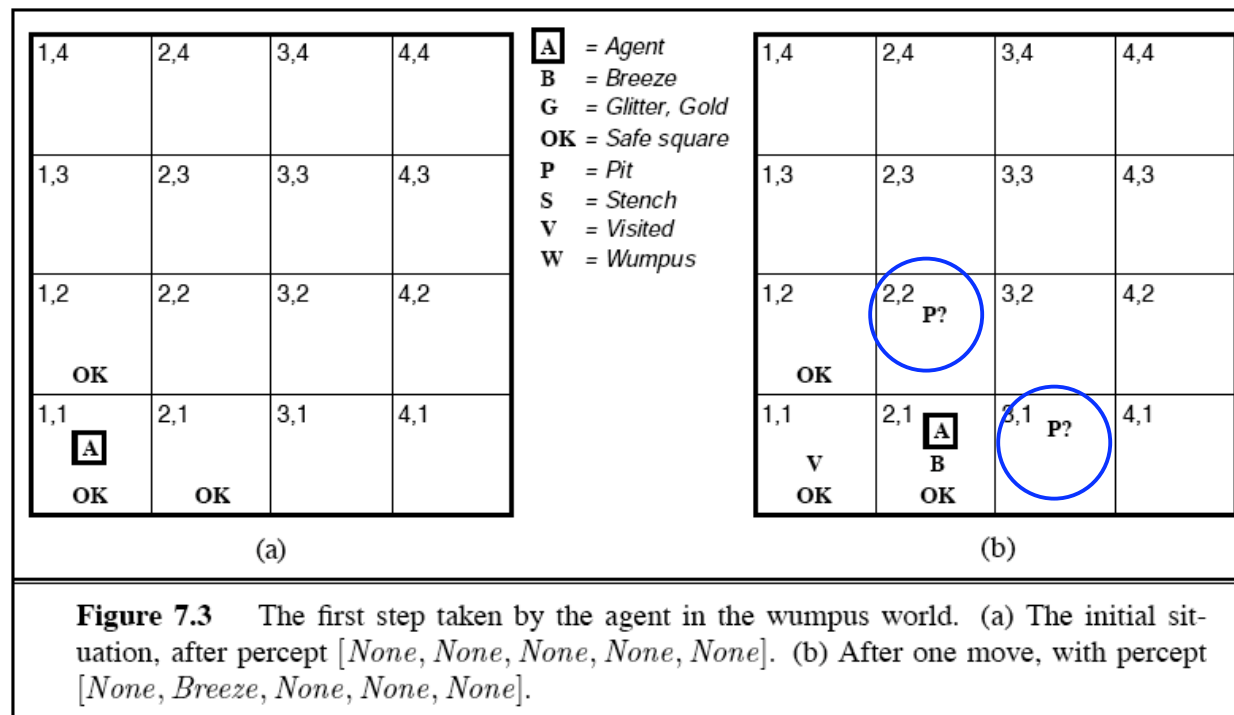[*Stench = {T,F}, Breeze = {T,F}, Glitter = {T,F}, Bumb = {T,F}, Scream = {T,F}*]

Example:
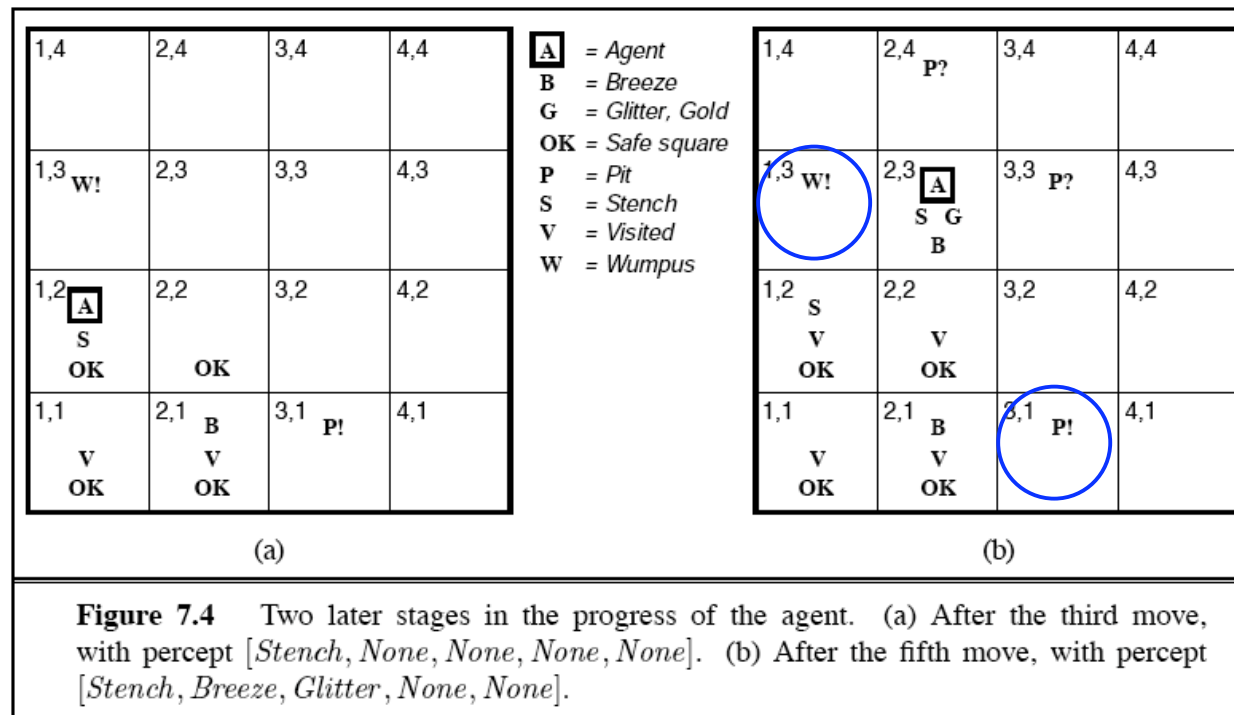
[*Stench, Breeze, None, None, None*]

# Logical Reasoning in the Wumpus World

1. In the beginning the agent only knows that it is in [1,1] and that [1,1] is safe

2. first percept [*None, None, None, None, None*] → neighboring squares are safe

3. careful agent decides to move to [2,1]

4. agent perceives a breeze (percept [*None, Breeze, None, None, None*]) → there must be a pit in either [2,2] or [3,1] or both



**Figure 7.3** The first step taken by the agent in the wumpus world. (a) The initial situation, after percept [*None, None, None, None, None*]. (b) After one move, with percept [*None, Breeze, None, None, None*].

5. only safe square which is known safe is [1,1] → prudent agent will turn around and go back and then proceed to [1,2]

6. in [1,2] new percept [*Stench, None, None, None, None*] → wumpus is nearby

7. wumpus cannot be in [1,1] and not in [2,2] either, otherwise stench in [2,1] → wumpus must be in [1,3]

8. lack of breeze in [1,2] → no pit in [2,2] → pit in [3,1]

9. …



**Figure 7.4**  Two later stages in the progress of the agent.  (a) After the third move, with percept [*Stench, None, None, None, None*].  (b) After the fifth move, with percept [*Stench, Breeze, Glitter, None, None*].

▸ This is fairly difficult inference, because it combines knowledge gained at different times in different places and relies on the lack of a percept to make one crucial step.

▸ The inference is beyond the abilities of most animals, but it is typical of the kind of reasoning that a logical agent does.

▸ In each case where the agent draws a conclusion from the available information, that conclusion is guaranteed to be correct if the available information is correct.

▸ Uncertainty arises in the wumpus world because the agent's sensors give only partial, local information about the world.

▸ Figure below shows a situation in which each of the three reachable squares [1,3], [2,2], and [3,1] might contain a pit.

▸ Often pure logical inference can conclude nothing about which square is most likely to be safe, so a logical agent might be forced to choose randomly.
We will see that probabilistic agent can do much better.

| 1,4 | 2,4 | 3,4 | 4,4 |
|-----|-----|-----|-----|
| 1,3 | 2,3 | 3,3 | 4,3 |
| 1,2 B  OK | 2,2 | 3,2 | 4,2 |
| 1,1  OK | 2,1 B  OK | 3,1 | 4,1 |

# Probabilistic Reasoning in the Wumpus World

▸ Calculate the probability that each of the three squares contains a pit.

▸ The relevant properties of the wumpus world are that
(1) a pit causes breezes in all neighboring squares, and
(2) each square other than [1,1] contains a pit with probability 0.2.

▸ <u>First step</u>: to identify the set of random variables we need:

– Boolean variable $P_{ij} = true$ iff square [i, j] contains pit
– Boolean variable $B_{ij} = true$ iff square [i, j] is breezy

include these variables in probability model only for the observed squares [1,1], [1,2], and [2,1].

▸ <u>Next step</u>: specify the full joint distribution, $\mathbf{P}(P_{1,1}, \dots, P_{4,4}, B_{1,1}, B_{1,2}, B_{2,1})$.

Applying the product rule, we have

$$\mathbf{P}(P_{1,1}, \dots, P_{4,4}, B_{1,1}, B_{1,2}, B_{2,1}) =$$

$$= \mathbf{P}(B_{1,1}, B_{1,2}, B_{2,1} \mid P_{1,1}, \dots, P_{4,4}) \, \mathbf{P}(P_{1,1}, \dots, P_{4,4})$$

$$\mathbf{P}(P_{1,1}, \ldots, P_{4,4}, B_{1,1}, B_{1,2}, B_{2,1}) =$$

$$= \mathbf{P}(B_{1,1}, B_{1,2}, B_{2,1} \mid P_{1,1}, \ldots, P_{4,4})\, \mathbf{P}(P_{1,1}, \ldots, P_{4,4})$$

▸ decomposition makes it very easy to see what joint probability values should be:
  − first term is conditional probability of a breeze configuration, given a pit configuration, equals 1 if the breezes are adjacent to the pits and 0 otherwise.
  − second term is the prior probability of a pit configuration

Each square contains a pit with probability 0.2, independently of the other squares; hence,

$$P(P_{1,1}, \ldots, P_{4,4}) = \prod_{i,j=1,1}^{4,4} P(P_{i,j})$$

For a configuration with $n$ pits, this is just

$$0.8 \times \ldots \times 0.8 \times \mathbf{0.2} \times 0.8 \times \ldots \times 0.8 \times \mathbf{0.2} \times 0.8 \times \ldots \times 0.8 \times \mathbf{0.2} \times 0.8 \times \ldots \times 0.8 =$$

$$= 0.2^3 \times 0.8^{16\text{-}3}$$

▸ in the situation in the figure the evidence consists of the observed breeze (or its absence) in each visited square combined with the fact that each such square contains no pit,

abbreviated as

$$b = \neg b_{1,1} \wedge b_{1,2} \wedge b_{2,1} \quad \text{and}$$
$$known = \neg p_{1,1} \wedge \neg p_{1,2} \wedge \neg p_{2,1}$$

▸ We are interested in answering queries such as **P(** $P_{1,3}$ | *known, b):*
how likely is it that [1,3] contains a pit, given the observations so far?

| 1,4 | 2,4 | 3,4 | 4,4 |
|-----|-----|-----|-----|
| 1,3 | 2,3 | 3,3 | 4,3 |
| 1,2 **B** OK | 2,2 | 3,2 | 4,2 |
| 1,1 OK | 2,1 **B** OK | 3,1 | 4,1 |

After finding a breeze in both [1,2] and [2,1] the agent is stuck - there is no safe place to explore

▸ To answer query, we can sum over entries from the full joint distribution.

Let *Unknown* be a composite variable consisting of the $P_{i,j}$ variables for squares other than the *Known* squares and the query square [1,3].

By Equation (1.6), we have

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

▸ The full joint probabilities have already been specified, so we are done that is, unless we care about computation.

There are 12 unknown squares; hence the summation contains $2^{12} = 4096$ terms.

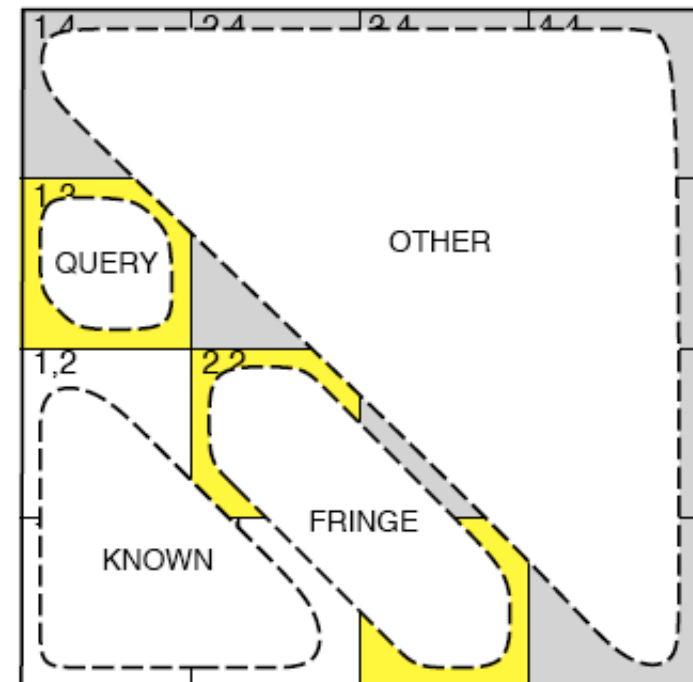In general, the summation grows exponentially with the number of squares.

▸ Are we missing something here? Are the other squares really relevant for our considerations?

▸ Not really! The contents of [4,4] don't affect whether [1,3] has a pit!

▸ Let *Fringe* be the variables (other than the query variable) that are adjacent to visited squares, in this case just [2,2] and [3,1].

▸ Also, let *Other* be the variables for the other unknown squares; there are 10 other squares, as shown in the figure

Define *Unknown* = *Fringe* ∪ *Other*

▸ The key insight is that the observed breezes are conditionally independent of the other variables

$$\mathbf{P}(\,b\,|\,P_{1,3}\,,\,Known,\,Unknown) =$$
$$= \mathbf{P}(\,b\,|\,P_{1,3}\,,\,Known,\,Fringe)$$

▸ To use the insight, we manipulate the query formula into a form in which the breezes are conditioned on all the other variables, and then we simplify using conditional independence

$$\mathbf{P}(P_{1,3} \mid known, b) =$$

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

$$= \alpha \sum_{fringe} \sum_{other} \mathbf{P}(b \mid P_{1,3}, known, fringe, other) \, \mathbf{P}(P_{1,3}, known, fringe, other)$$

conditional independence

$$= \alpha \sum_{fringe} \sum_{other} \mathbf{P}(b \mid P_{1,3}, known, fringe) \, \mathbf{P}(P_{1,3}, known, fringe, other)$$

first term independent of other variables → move summation inwards

$$= \alpha \sum_{fringe} \mathbf{P}(b \mid P_{1,3}, known, fringe) \sum_{other} \mathbf{P}(P_{1,3}, known, fringe, other)$$

independence → factoring out prior, reordering terms

$$= \alpha \sum_{fringe} \mathbf{P}(b \mid P_{1,3}, known, fringe) \sum_{other} \mathbf{P}(P_{1,3}) P(known) P(fringe) P(other)$$

$$\mathbf{P}(P_{1,3} \mid known, b) =$$

$$= \ldots$$

$$= \alpha \sum_{fringe} \mathbf{P}(b \mid P_{1,3}, known, fringe) \sum_{other} \mathbf{P}(P_{1,3}) P(known) P(fringe) P(other)$$

$$= \alpha \, P(known) \mathbf{P}(P_{1,3}) \sum_{fringe} \mathbf{P}(b \mid P_{1,3}, known, fringe) \, P(fringe) \sum_{other} P(other)$$

$$= \alpha' \mathbf{P}(P_{1,3}) \sum_{fringe} \mathbf{P}(b \mid P_{1,3}, known, fringe) \, P(fringe)$$

$$\mathbf{P}(P_{1,3} \mid known, b) =$$

$$= \dots$$

$$= \alpha \sum_{fringe} \mathbf{P}(b \mid P_{1,3}, known, fringe) \sum_{other} \mathbf{P}(P_{1,3}) P(known) P(fringe) P(other)$$

$$= \alpha \ P(known) \mathbf{P}(P_{1,3}) \sum_{fringe} \mathbf{P}(b \mid P_{1,3}, known, fringe) \ P(fringe) \sum_{other} P(other)$$

$$= \alpha' \mathbf{P}(P_{1,3}) \sum_{fringe} \mathbf{P}(b \mid P_{1,3}, known, fringe) \ P(fringe)$$

?

?

$$\mathbf{P}(P_{1,3} \mid known, b) =$$

$$= \ldots$$

$$= \alpha \sum_{fringe} \mathbf{P}(b \mid P_{1,3}, known, fringe) \sum_{other} \mathbf{P}(P_{1,3}) P(known) P(fringe) P(other)$$

$$= \alpha \; P(known) \mathbf{P}(P_{1,3}) \sum_{fringe} \mathbf{P}(b \mid P_{1,3}, known, fringe) \; P(fringe) \sum_{other} P(other)$$

$$= \alpha' \mathbf{P}(P_{1,3}) \sum_{fringe} \mathbf{P}(b \mid P_{1,3}, known, fringe) \; P(fringe)$$

probability of observing $b$ given the pit configuration

$other = ?$

$P(other) = ?$

$$\mathbf{P}(P_{1,3} \mid known, b) =$$

$$= \ldots$$

$$= \alpha \sum_{fringe} \mathbf{P}(b \mid P_{1,3}, known, fringe) \sum_{other} \mathbf{P}(P_{1,3}) P(known) P(fringe) P(other)$$

$$= \alpha \, P(known) \mathbf{P}(P_{1,3}) \sum_{fringe} \mathbf{P}(b \mid P_{1,3}, known, fringe) \, P(fringe) \sum_{other} P(other) \leftarrow$$

$$= \alpha' \mathbf{P}(P_{1,3}) \sum_{fringe} \mathbf{P}(b \mid P_{1,3}, known, fringe) \, P(fringe)$$
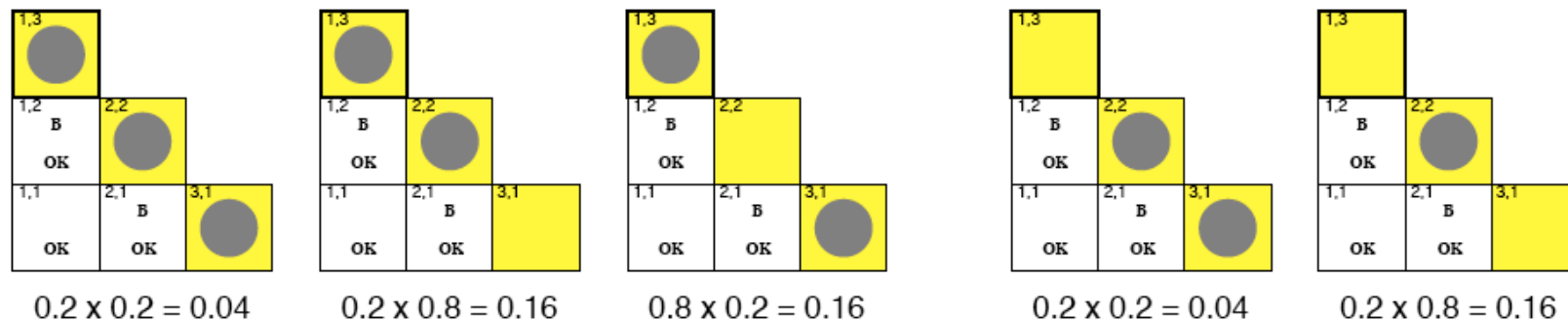
$$\overset{?}{=} 1$$

probability of observing $b$ given the pit configuration

$$= \alpha' \mathbf{P}(P_{1,3})(\mathbf{P}(b \mid P_{1,3}, known, p_{2,2}, p_{3,1}) P(p_{2,2} \wedge p_{3,1}) +$$
$$\mathbf{P}(b \mid P_{1,3}, known, p_{2,2}, \neg p_{3,1}) P(p_{2,2} \wedge \neg p_{3,1}) +$$
$$\mathbf{P}(b \mid P_{1,3}, known, \neg p_{2,2}, p_{3,1}) P(\neg p_{2,2} \wedge p_{3,1}) +$$
$$\mathbf{P}(b \mid P_{1,3}, known, \neg p_{2,2}, \neg p_{3,1}) P(\neg p_{2,2} \wedge \neg p_{3,1}))$$

$$= \alpha' \left\langle \begin{array}{l} 0.2(1 \times (0.2 \times 0.2) + 1 \times (0.2 \times 0.8) + 1 \times (0.8 \times 0.2) + 0 \times (0.8 \times 0.8)), \\ 0.8(1 \times (0.2 \times 0.2) + 1 \times (0.2 \times 0.8) + 0 \times (0.8 \times 0.2) + 0 \times (0.8 \times 0.8)) \end{array} \right\rangle$$

▸ there are just four terms in the summation over the fringe variables $P_{2,2}$ and $P_{3,1}$ . Use of independence and conditional independence has completely eliminated the other squares from consideration.



0.2 x 0.2 = 0.04          0.2 x 0.8 = 0.16          0.8 x 0.2 = 0.16          0.2 x 0.2 = 0.04          0.2 x 0.8 = 0.16

Consistent models for the fringe variables $P_{2,2}$ and $P_{3,1}$, showing $P(\text{fringe})$ for each model: (a) three models with $P_{3,1} = \text{true}$ showing two or three pits and (a) two models $P_{3,1} = \text{false}$ showing one or two pits.

$$\mathbf{P}(P_{1,3}|known, b) = \alpha' \langle 0.2(0.04 + 0.16 + 0.16), \ 0.8(0.04 + 0.16) \rangle$$
$$\approx \langle 0.31, 0.69 \rangle$$

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

# 1.8  Summary

▸ Uncertainty arises because of both laziness and ignorance. It is inescapable in complex, dynamic, or inaccessible worlds.

▸ Uncertainty means that many of the simplifications that are possible with deductive inference are no longer valid.

▸ Probabilities express the agent's inability to reach a definite decision regarding the truth of a sentence. Probabilities summarize the agent's beliefs.

▸ Basic probability statements include prior probabilities and conditional probabilities over simple and complex propositions.

▸ The full joint probability distribution specifies the probability of each complete assignment of values to random variables. It is usually too large to create or use in its explicit form.

▸ The axioms of probability constrain the possible assignments of probabilities to propositions. An agent that violates the axioms will behave irrationally in some circumstances.

‣ When the full joint distribution is available, it can be used to answer queries simply by adding up entries for the atomic events corresponding to the query propositions.

‣ Absolute independence between subsets of random variables might allow the full joint distribution to be factored into smaller joint distributions. This could greatly reduce complexity, but seldom occurs in practice.

‣ Bayes' rule allows unknown probabilities to be computed from known conditional probabilities, usually in the causal direction. Applying Bayes' rule with many pieces of evidence will in general run into the same scaling problems as does the full joint distribution.

‣ Conditional independence brought about by direct causal relationships in the domain might allow the full joint distribution to be factored into smaller, conditional distributions.

‣ A wumpus-world agent can calculate probabilities for unobserved aspects of the world and use them to make better decisions than a purely logical agent makes.