

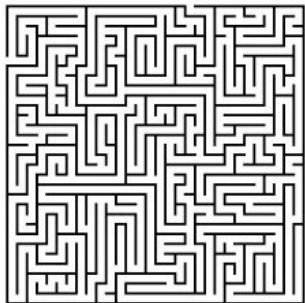
# Uncertainty and Probability

# Lecture outline

- Uncertainty
- Probability
  - Combinatorics
  - Joint probability
  - Marginal probability
  - Conditional probability
  - Independence
  - Conditional independence
  - Bayes' rule

# Uncertainty

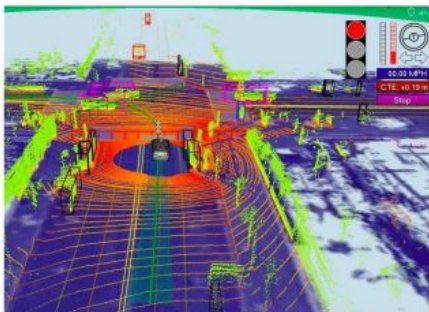
# Deterministic Reasoning vs Probabilistic Reasoning



8			4	6		7
	1				4	
5	9		3		6	5
			7			
	4	8	2		1	3
	5	2				9
3		1				
		9		2		5



VS.



# Uncertainty

- If I take my umbrella, I won't get wet in the rain
  - **Deterministic reasoning**
- But there is a chance that:
  - The umbrella might brake and I'll get wet
  - There will be strong winds spraying water everywhere and I'll get wet
  - **Probabilistic reasoning**
- Agents almost never have perfect information about the world
  - In those situations, the agent must reason under **uncertainty**
  - Uncertainty can also arise because of the agent's incorrect/incomplete understanding of its environment

# Making decisions under uncertainty

Let **Action**  $A_t$  = *"leave for airport  $t$  minutes before the flight"*

Will  $A_t$  get me to the flight in time?

## Problems:

- Partial observability (road state, other drivers,...)
- Noisy sensors (unreliable traffic reports, weather reports,...)
- Uncertainty in action outcomes (flat tyre, mechanical failure,...)
- Complexity of modeling and predicting traffic

# Making decisions under uncertainty

A purely logical approach either

1. Risks falsehood
  - a.  $A_{25}$  will get me there on time, or
2. Leads to conclusions that are too weak for decision making
  - a.  $A_{25}$  will get me there on time if there is no accident on the way and there is no rain and there is no problem with the car, and ...

$A_{2500}$  might get me there in time for the flight but I might have to stay at the airport for far too long!

# Reasoning under uncertainty

A **rational agent** is one that makes rational decisions — in order to maximize its performance measure)

A rational decision depends on:

- the **relative importance** of various goals
- the **likelihood** they will be achieved
- the **degree** to which they will be achieved



# Handling uncertain knowledge

Reasons [First Order Logic](#) based approaches fail to cope with domains like, for instance, medical diagnosis:

- **Laziness:** too much work to write complete axioms, or too hard to work with the enormous sentences that result
- **Theoretical Ignorance:** The available knowledge of the domain is incomplete
- **Practical Ignorance:** The theoretical knowledge of the domain is complete but some evidential facts are missing

# Degrees of belief

- In several real-world domains the agent's knowledge can only provide a **degree of belief** in the relevant sentences
  - The agent cannot say whether a sentence is true, but only that it is true  $x\%$  of the time
- The main tool for handling degrees of belief is **Probability Theory**
- The use of probability summarizes the uncertainty that stems from our laziness or ignorance about the domain

# Methods for handling uncertainty

- Default or non-monotonic logic:
  - E.g., assume  $A_{25}$  works unless contradicted by evidence
  - Issues: What assumptions are reasonable? How to handle contradiction?
- Probability
  - E.g., Given the available evidence,  $A_{25}$  will get me there on time with **probability 0.04**

# Making decisions under uncertainty

- Suppose the agent believes the following:
  - $P(A_{25} \text{ gets me there on time}) = 0.04$
  - $P(A_{90} \text{ gets me there on time}) = 0.70$
  - $P(A_{120} \text{ gets me there on time}) = 0.95$
  - $P(A_{1440} \text{ gets me there on time}) = 0.9999$
- Which action should the agent choose?
  - Depends on preferences for missing flight vs. time spent waiting
  - Encapsulated by a **utility function**
    - Attempts to quantify satisfaction/happiness
- The agent should choose the action that maximizes the **expected utility**:
  - $P(A_t \text{ succeeds}) * U(A_t \text{ succeeds}) + P(A_t \text{ fails}) * U(A_t \text{ fails})$

# Making decisions under uncertainty


- More generally: the ***expected utility of an action*** is defined as:
- $EU(a) = \sum_{\text{outcomes of } a} P(\text{outcome} \mid a) * U(\text{outcome})$
- Utility theory is used to represent and infer preferences
- Decision theory = probability theory + utility theory

# Probability

# Probability

- Probabilistic assertions summarize effects of
  - laziness: failure to enumerate exceptions, etc.
  - ignorance: lack of relevant facts, initial conditions, etc.
- Subjective or Bayesian probability:
  - Probabilities relate propositions to one's own state of knowledge
  - e.g.,  $P(A_{25} \text{ gets me there on time} \mid \text{no reported accidents}) = 0.06$
- Probabilities of propositions change with new evidence:
  - E.g.,  $P(A_{25} \text{ gets me there on time} \mid \text{no reported accidents, 5 a.m.}) = 0.15$

# Probability basics

- Begin with a set  $\Omega$  – the sample space
  - e.g., 6 possible rolls of a die.
- $\omega \in \Omega$  is a **sample point/possible world/atomic event**
- Atomic Event
  - An atomic event is a complete specification of the state of the world.
  - E.g. if the world is only concerned about two Boolean variables A and B. There are four possible atomic events:  $A \wedge B$ ,  $\neg A \wedge B$ ,  $A \wedge \neg B$  and  $\neg A \wedge \neg B$  
- E.g. for the roll of the die, 1, 2 3,..., 6 are all atomic events
- Atomic events are mutually exclusive and collectively exhaustive.
  - $P(\omega_i, \omega_j) = 0$  for all  $i \neq j$
  - $\sum_{(\omega \in \Omega)} P(\omega) = 1$



# Probability basics (2)

- A **probability space** or **probability model** is a sample space with an assignment  $P(\omega)$  for every  $\omega \in \Omega$  s.t.
  - $0 \leq P(\omega) \leq 1$
  - $\sum_{\omega} P(\omega) = 1$
  - e.g.,  $P(1)=P(2)=P(3)=P(4)=P(5)=P(6)=1/6$ .
- An **event**  $A$  is any subset of  $\Omega$ 
  - $P(A) = \sum_{(\omega \in A)} P(\omega)$
  - E.g.,  $P(\text{die roll} < 4) = P(1) + P(2) + P(3) = 1/6 + 1/6 + 1/6 = 1/2$

# Random variables

- A random variable is a variable whose possible values (Domain ) are numerical outcomes of a random phenomenon.
  - E.g.  $X$  = outcome of a roll of a dice.
  - $\text{Domain}(X) = \{1, 2, 3, 4, 5, 6\}$
- A Random variable can be boolean, discrete or continuous, depending on its domain.
- **$P$**  induces a probability distribution for any r.v.  $X$ :
  - e.g.,  $P(X=1) = 1/6$ ,  $P(X=2) = 1/6$ , ...,  $P(X=6) = 1/6$

# Propositions

- Think of a proposition as the event (set of sample points) where the proposition is true
  - Given the random variable,  $X$  = outcome of a roll of a dice
- The proposition,  $X\_is\_Odd$ , can be thought of as the event,
  - $X\_is\_Odd = \{ \omega \in \Omega \mid X \text{ is Odd} \}$
- Proposition = disjunction of atomic events in which it is true. E.g.
  - $X\_is\_Odd \Leftrightarrow (X=1) \vee (X=3) \vee (X=5)$
  - $(a \vee b) \Leftrightarrow (\neg a \sqcap b) \vee (a \sqcap \neg b) \vee (a \sqcap b)$

# Syntax of Propositions

- **Propositional** or **Boolean** random variables
  - e.g., Cavity (do I have a cavity?)
  - **Cavity =true** is a proposition, also written *cavity*
- **Discrete** random variables (finite or infinite)
  - e.g., **Weather** is one of **<sunny, rain, cloudy, snow>**
  - **Weather = rain** is a proposition
  - Values must be exhaustive and mutually exclusive
- **Continuous** random variables (bounded or unbounded)
  - e.g., **Temp=21.6**; also allow, e.g., **Temp < 22.0**.
- Arbitrary Boolean combinations of basic propositions

# 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:
  - $P(\text{Weather} = \text{sunny}) = 0.7$
  - $P(\text{Weather} = \text{rain}) = 0.2$
  - $P(\text{Weather} = \text{cloudy}) = 0.08$
  - $P(\text{Weather} = \text{snow}) = 0.02$
- **$P(\text{Weather}) = \langle 0.7, 0.2, 0.08, 0.02 \rangle$**  (normalized, i.e., **sums to 1**)

## Prior probability (2)

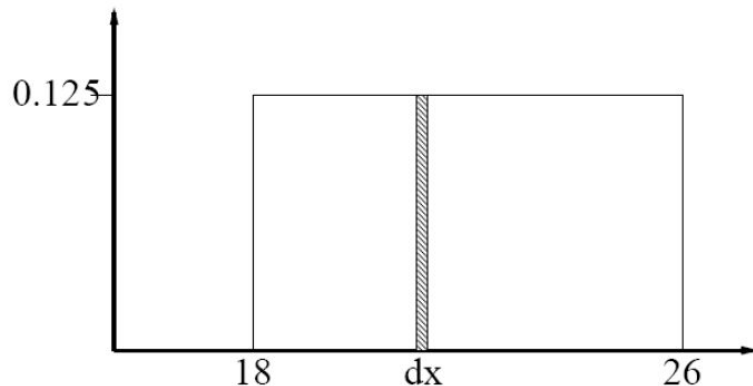
- Joint probability distribution for a set of r.v.s gives the probability of every atomic event on those r.v.s (i.e., every sample point)
- $P(\text{Weather}, \text{Cavity})$  is a  $4 \times 2$  matrix of values:

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

- Every question about a domain can be answered by the joint distribution because every event is a sum of sample points

# Probability for continuous variables

- Probability distributions of continuous r.v s are specified by **probability density functions (pdfs)**
  - E.g. Uniform distribution between 18 and 26



Here  $P$  is a **density**; integrates to 1.

$P(X = 20.5) = 0.125$  really means

$$\lim_{dx \rightarrow 0} P(20.5 \leq X \leq 20.5 + dx)/dx = 0.125$$

# Joint probability distributions

- A joint distribution is an assignment of probabilities to every possible atomic event

Atomic event	P
$Cavity = false \wedge Toothache = false$	0.8
$Cavity = false \wedge Toothache = true$	0.1
$Cavity = true \wedge Toothache = false$	0.05
$Cavity = true \wedge Toothache = true$	0.05

- Suppose we have a joint distribution of  $n$  random variables with domain sizes  $d$ 
  - What is the size of the probability table?
  - Impossible to write out completely for all but the smallest distributions



# Notation

- $P(X_1 = x_1, X_2 = x_2, \dots, X_n = x_n)$  refers to a single entry (atomic event) in the joint probability distribution table
  - Shorthand:  $P(x_1, x_2, \dots, x_n)$
- $P(X_1, X_2, \dots, X_n)$  refers to the entire joint probability distribution table
- $P(A)$  can also refer to the probability of an event
  - E.g.,  $X_1 = x_1$  is an event

# Marginal probability distributions

- From the **joint** distribution **P(X,Y)** we can find the **marginal** distributions **P(X)** and **P(Y)**
- To find  $P(X = x)$ , sum the probabilities of all atomic events where  $X = x$ :

$$\begin{aligned} P(X = x) &= P((X = x \wedge Y = y_1) \vee \dots \vee (X = x \wedge Y = y_n)) \\ &= P((x, y_1) \vee \dots \vee (x, y_n)) = \sum_{i=1}^n P(x, y_i) \end{aligned}$$

- This is called **marginalization** (we are marginalizing out all the variables except X)

# Marginal probability distributions

From the **joint** distribution **P(X,Y)** we can find the **marginal** distributions **P(X)** and **P(Y)**

<b>P(Cavity, Toothache)</b>	
<i>Cavity = false <math>\wedge</math> Toothache = false</i>	0.8
<i>Cavity = false <math>\wedge</math> Toothache = true</i>	0.1
<i>Cavity = true <math>\wedge</math> Toothache = false</i>	0.05
<i>Cavity = true <math>\wedge</math> Toothache = true</i>	0.05

<b>P(Cavity)</b>	
<i>Cavity = false</i>	?
<i>Cavity = true</i>	?

<b>P(Toothache)</b>	
<i>Toothache = false</i>	?
<i>Toothache = true</i>	?

# 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
- If we know more, e.g., there is a rough spot on the tooth, then we have
  - $P(\text{cavity} \mid \text{toothache}, \text{rough\_spot}) = 0.9$
- New evidence may be irrelevant, allowing simplification, e.g.,
  - $P(\text{cavity} \mid \text{toothache}, \text{sunny}) = P(\text{cavity} \mid \text{toothache}) = 0.8$
- This kind of inference, sanctioned by domain knowledge, is crucial

## Conditional probability (2)

- Definition of conditional probability: 
$$P(a/b) = \frac{P(a \wedge b)}{P(b)} \quad \text{if } P(b) > 0$$
- **Product rule** gives an alternative formulation:
  - $P(a \wedge b) = P(a | b) P(b) = P(b | a) P(a)$
- A general version holds for whole distributions, e.g.,
  - $P(\text{Weather}, \text{Cavity}) = P(\text{Weather} | \text{Cavity}) P(\text{Cavity})$
  - $P(A, B) = P(A|B).P(B) = P(B|A).P(A)$

# Chain rule

Using successive application of the product rule

- Product rule

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

- Chain rule

$$\begin{aligned} \mathbf{P}(X_1, \dots, X_n) &= \mathbf{P}(X_n | X_{n-1}, \dots, X_1) \mathbf{P}(X_{n-1}, \dots, X_1) \\ &= \mathbf{P}(X_n | X_{n-1}, \dots, X_1) \mathbf{P}(X_{n-1} | X_{n-2}, \dots, X_1) \mathbf{P}(X_{n-2}, \dots, X_1) \\ &= \mathbf{P}(X_n | X_{n-1}, \dots, X_1) \mathbf{P}(X_{n-1} | X_{n-2}, \dots, X_1) \dots \mathbf{P}(X_2 | X_1) \mathbf{P}(X_1) \\ &= \\ &\quad \prod_{i=1}^n P(X_i | X_{i-1}, \dots, X_1) \end{aligned}$$

# 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>	<b>.108</b>	<b>.012</b>	<b>.072</b>	<b>.008</b>
$\neg$ <i>cavity</i>	<b>.016</b>	<b>.064</b>	<b>.144</b>	<b>.576</b>

- Catch means the dentists steel probe getting caught in the tooth.
- For any **proposition  $\phi$** , sum the **atomic events** where it is true:
  - $P(\phi) = \sum_{\omega: \omega \models \phi} P(\omega)$
  - $P(\text{toothache}) = ?$
  - $P(\text{cavity} \vee \text{toothache}) = ?$
  - $P(\sim \text{cavity} \mid \text{toothache}) = ?$

## Inference by enumeration (2)

- 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>	<b>.108</b>	<b>.012</b>	<b>.072</b>	<b>.008</b>
$\neg$ <i>cavity</i>	<b>.016</b>	<b>.064</b>	<b>.144</b>	<b>.576</b>

- For any proposition  $\phi$ , 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 (3)

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

For any proposition  $\phi$ , sum the atomic events where it is true:

$$P(\phi) = \sum_{\omega: \omega \models \phi} P(\omega)$$

$$P(\text{cavity} \vee \text{toothache}) = 0.108 + 0.012 + 0.072 + 0.008 + 0.016 + 0.064 = 0.28$$

# Inference by enumeration (4)

- 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>	<b>.108</b>	<b>.012</b>	<b>.072</b>	<b>.008</b>
$\neg$ <i>cavity</i>	<b>.016</b>	<b>.064</b>	<b>.144</b>	<b>.576</b>

- Can also compute conditional probabilities:

$$\begin{aligned} P(\neg \text{cavity} | \text{toothache}) &= \frac{P(\neg \text{cavity} \wedge \text{toothache})}{P(\text{toothache})} \\ &= \frac{0.016 + 0.064}{0.108 + 0.012 + 0.016 + 0.064} \\ &= 0.4 \\ P(\text{cavity} | \text{toothache}) &= \frac{P(\text{cavity} \wedge \text{toothache})}{P(\text{toothache})} = 0.6 \end{aligned}$$

# 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>	<b>.108</b>	<b>.012</b>	<b>.072</b>	<b>.008</b>
$\neg$ <i>cavity</i>	<b>.016</b>	<b>.064</b>	<b>.144</b>	<b>.576</b>

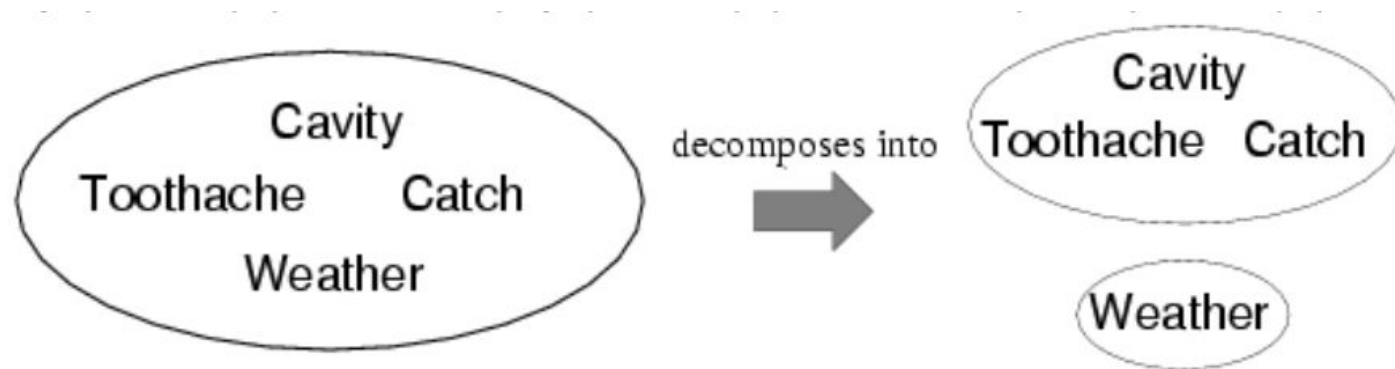
- Denominator can be viewed as a **normalization constant**  $\alpha$ 
  - $P(\text{Cavity} \mid \text{toothache}) = \alpha P(\text{Cavity}, \text{toothache})$
  - $= \alpha [P(\text{Cavity}, \text{toothache}, \text{catch}) + P(\text{Cavity}, \text{toothache}, \sim \text{catch})]$
  - $= \alpha [0.108 + 0.012]$
  - $= \alpha 0.12$
- General idea: compute distribution on query variable by fixing **evidence variables** and summing over **hidden variables**

# Inference by enumeration (5)

- Let **X** be the set of all the variables.
  - Typically, we want 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:
  - Worst-case time complexity  $O(2^n)$

# Independence

- A and B are independent iff:
  - $P(A|B) = P(A)$  or  $P(B|A) = P(B)$  or  $P(A, B) = P(A) P(B)$



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

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

# Conditional independence

- $P(\text{Toothache}, \text{Cavity}, \text{Catch})$  has  $2^3 = 8$  independent entries
- If someone has a cavity, the probability that a probe catches in it doesn't depend on whether he has a toothache:
  - $P(\text{catch} \mid \text{toothache}, \text{cavity}) = P(\text{catch} \mid \text{cavity})$
- The same independence holds if he hasn't got a cavity:
  - $P(\text{catch} \mid \text{toothache}, \sim\text{cavity}) = P(\text{catch} \mid \sim\text{cavity})$
- ***Catch*** is **conditionally independent** of ***Toothache*** given ***Cavity***:
- $P(\text{Catch} \mid \text{Toothache}, \text{Cavity}) = P(\text{Catch} \mid \text{Cavity})$
- Equivalent statements:
  - $P(\text{Toothache} \mid \text{Catch}, \text{Cavity}) = P(\text{Toothache} \mid \text{Cavity})$
  - $P(\text{Toothache}, \text{Catch} \mid \text{Cavity}) = P(\text{Toothache} \mid \text{Cavity}) P(\text{Catch} \mid \text{Cavity})$

# Conditional independence (2)

- Write out full joint distribution using chain rule:
- $P(\text{Toothache}, \text{Catch}, \text{Cavity})$ 
  - $= P(\text{Toothache} \mid \text{Catch}, \text{Cavity}) P(\text{Catch}, \text{Cavity})$
  - $= P(\text{Toothache} \mid \text{Catch}, \text{Cavity}) P(\text{Catch} \mid \text{Cavity}) P(\text{Cavity})$
  - $= P(\textbf{Toothache} \mid \textbf{Cavity}) P(\textbf{Catch} \mid \textbf{Cavity}) P(\textbf{Cavity})$
- I.e.,  $\mathbf{2} + \mathbf{2} + \mathbf{1} = \mathbf{5}$  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.

# Bayes' Rule

**Product rule**  $P(a \wedge b) = P(a \mid b) P(b) = P(b \mid a) P(a)$

**Bayes' rule:**

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

or in **distribution form:**

$$P(Y \mid X) = \frac{P(X \mid Y)P(Y)}{P(X)}$$

Useful for assessing diagnostic probability from causal probability:

$$P(Cause \mid Effect) = \frac{P(Effect \mid Cause)P(Cause)}{P(Effect)}$$

E.g., let M be meningitis, S be stiff neck:

$$P(m|s) = P(s|m) P(m) / P(s) = 0.8 \times 0.0001 / 0.1 = 0.0008$$





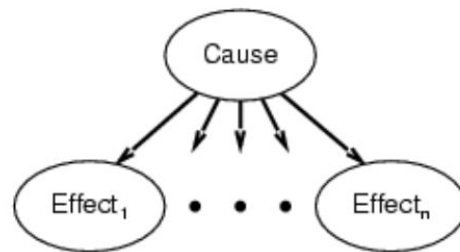
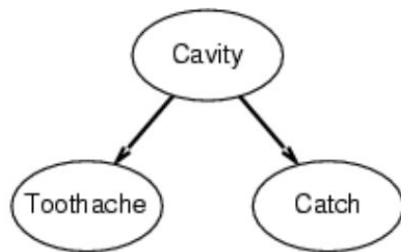
# Bayes' rule and conditional independence

$P(\text{Cavity} \mid \text{toothache} \wedge \text{catch})$

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

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

$$P(\text{Cause}, \text{Effect}_1, \dots, \text{Effect}_n) = P(\text{Cause}) \prod_i P(\text{Effect}_i \mid \text{Cause})$$



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