

Bayesian Reasoning

Chapters 12 & 13



[Thomas Bayes, 1701-1761](#)

Today's topics

- Motivation
- Review probability theory
- Bayesian inference
 - From the joint distribution
 - Using independence/factoring
 - From sources of evidence
- Naïve Bayes algorithm for inference and classification tasks

Motivation: causal reasoning



- As the sun rises, the rooster crows
 - Does this correlation imply causality?
 - If so, which way does it go?
- The evidence can come from
 - Probabilities and Bayesian reasoning
 - Common sense knowledge
 - Experiments
- Bayesian Belief Networks (BBNs) are useful for modeling causal reasoning

Many Sources of Uncertainty

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

Decision making with uncertainty

Rational behavior: for each possible action:

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

Consider

- Your house has an alarm system
- It should go off if a burglar breaks into the house
- It can go off if there is an earthquake
- How can we predict what's happened if the alarm goes off?
 - Someone has broken in!
 - It's a minor earthquake



Probability theory 101

- **Random variables:**

- Domain

- **Atomic event:**

- complete specification of state

- **Prior probability:**

- degree of belief without any other evidence or info

- **Joint probability:**

- matrix of combined probabilities of set of variables

- Alarm, Burglary, Earthquake

- Boolean (these) or discrete (0-9), continuous (float)

- Alarm=T \wedge Burglary=T \wedge Earthquake=F
alarm \wedge burglary \wedge \neg earthquake

- P(Burglary) = 0.1

- P(Alarm) = 0.1

- P(earthquake) = 0.000003

- P(Alarm, Burglary) =

	alarm	\neg alarm
burglary	.09	.01
\neg burglary	.1	.8

Probability theory 101

	alarm	¬alarm
burglary	.09	.01
¬burglary	.1	.8

- **Conditional probability:** prob. of effect given causes
- **Computing conditional probs:**
 - $P(a | b) = P(a \wedge b) / P(b)$
 - $P(b)$: **normalizing** constant
- **Product rule:**
 - $P(a \wedge b) = P(a | b) * P(b)$
- **Marginalizing:**
 - $P(B) = \sum_a P(B, a)$
 - $P(B) = \sum_a P(B | a) P(a)$ (**conditioning**)
- $P(\text{burglary} | \text{alarm}) = .47$
 $P(\text{alarm} | \text{burglary}) = .9$
- $P(\text{burglary} | \text{alarm}) = P(\text{burglary} \wedge \text{alarm}) / P(\text{alarm}) = .09 / .19 = .47$
- $P(\text{burglary} \wedge \text{alarm}) = P(\text{burglary} | \text{alarm}) * P(\text{alarm}) = .47 * .19 = .09$
- $P(\text{alarm}) = P(\text{alarm} \wedge \text{burglary}) + P(\text{alarm} \wedge \neg\text{burglary}) = .09 + .1 = .19$

Probability theory 101

	alarm	-alarm
burglary	.09	.01
-burglary	.1	.8

- **Conditional probability:** prob. of effect given causes

- **Computing conditional probs:**

- $P(a | b) = P(a \wedge b) / P(b)$

- $P(b)$: **normalizing** constant

- **Product rule:**

- $P(a \wedge b) = P(a | b) * P(b)$

- **Marginalizing:**

- $P(B) = \sum_a P(B, a)$

- $P(B) = \sum_a P(B | a) P(a)$
(**conditioning**)

- $P(\text{burglary} | \text{alarm}) = .47$
 $P(\text{alarm} | \text{burglary}) = .9$

- $P(\text{burglary} | \text{alarm}) =$
 $P(\text{burglary} \wedge \text{alarm}) / P(\text{alarm})$
 $= .09 / .19 = .47$

- $P(\text{burglary} \wedge \text{alarm}) =$
 $P(\text{burglary} | \text{alarm}) * P(\text{alarm})$
 $= .47 * .19 = .09$

- $P(\text{alarm}) =$
 $P(\text{alarm} \wedge \text{burglary}) +$
 $P(\text{alarm} \wedge \neg \text{burglary})$
 $= .09 + .1 = .19$

Example: Inference from the joint

	alarm		¬alarm	
	earthquake	¬earthquake	earthquake	¬earthquake
burglary	.01	.08	.001	.009
¬burglary	.01	.09	.01	.79

$$\begin{aligned} P(\text{burglary} \mid \text{alarm}) &= \alpha P(\text{burglary}, \text{alarm}) \\ &= \alpha [P(\text{burglary}, \text{alarm}, \text{earthquake}) + P(\text{burglary}, \text{alarm}, \neg\text{earthquake})] \\ &= \alpha [(.01, .01) + (.08, .09)] \\ &= \alpha [(.09, .1)] \end{aligned}$$

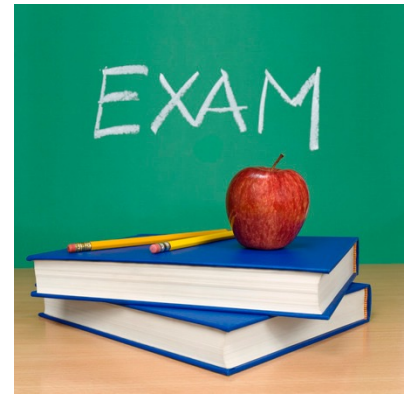
Since $P(\text{burglary} \mid \text{alarm}) + P(\neg\text{burglary} \mid \text{alarm}) = 1$, $\alpha = 1/(\text{.09} + \text{.1}) = 5.26$
(i.e., $P(\text{alarm}) = 1/\alpha = \text{.19}$ – **quizlet**: how can you verify this?)

$$P(\text{burglary} \mid \text{alarm}) = \text{.09} * 5.26 = \text{.474}$$

$$P(\neg\text{burglary} \mid \text{alarm}) = \text{.1} * 5.26 = \text{.526}$$

Consider

- A student has to take an exam
 - She might be smart
 - She might have studied
 - She may be prepared for the exam
- How are these related?
- We can collect joint probabilities for the three events
 - Measure prepared as “got a passing grade”



Exercise:

Inference from the joint



$p(\text{smart} \wedge \text{study} \wedge \text{prepared})$	smart		\neg smart	
	study	\neg study	study	\neg study
prepared	.432	.16	.084	.008
\neg prepared	.048	.16	.036	.072

Each of the eight highlighted boxes has the joint probability for the three values of smart, study, prepared

Queries:

- What is the prior probability of *smart*?
- What is the prior probability of *study*?
- What is the conditional probability of *prepared*, given *study* and *smart*?

Exercise:

Inference from the joint



$p(\text{smart} \wedge \text{study} \wedge \text{prepared})$	smart		$\neg\text{smart}$	
	study	$\neg\text{study}$	study	$\neg\text{study}$
prepared	.432	.16	.084	.008
$\neg\text{prepared}$.048	.16	.036	.072

Queries:

- What is the prior probability of *smart*?
- What is the prior probability of *study*?
- What is the conditional probability of *prepared*, given *study* and *smart*?

$$p(\text{smart}) = .432 + .16 + .048 + .16 = \mathbf{0.8}$$

Exercise:

Inference from the joint



$p(\text{smart} \wedge \text{study} \wedge \text{prepared})$	smart		\neg smart	
	study	\neg study	study	\neg study
prepared	.432	.16	.084	.008
\neg prepared	.048	.16	.036	.072

Queries:

- What is the prior probability of *smart*?
- **What is the prior probability of *study*?**
- What is the conditional probability of *prepared*, given *study* and *smart*?

Exercise:

Inference from the joint



$p(\text{smart} \wedge \text{study} \wedge \text{prepared})$	smart		\neg smart	
	study	\neg study	study	\neg study
prepared	.432	.16	.084	.008
\neg prepared	.048	.16	.036	.072

Queries:

- What is the prior probability of *smart*?
- **What is the prior probability of *study*?**
- What is the conditional probability of *prepared*, given *study* and *smart*?

$$p(\text{study}) = .432 + .048 + .084 + .036 = 0.6$$

Exercise:

Inference from the joint



$p(\text{smart} \wedge \text{study} \wedge \text{prepared})$	smart		\neg smart	
	study	\neg study	study	\neg study
prepared	.432	.16	.084	.008
\neg prepared	.048	.16	.036	.072

Queries:

- What is the prior probability of *smart*?
- What is the prior probability of *study*?
- **What is the conditional probability of *prepared*, given *study* and *smart*?**

Exercise:

Inference from the joint



$p(\text{smart} \wedge \text{study} \wedge \text{prepared})$	smart		$\neg\text{smart}$	
	study	$\neg\text{study}$	study	$\neg\text{study}$
prepared	.432	.16	.084	.008
$\neg\text{prepared}$.048	.16	.036	.072

Queries:

- What is the prior probability of *smart*?
- What is the prior probability of *study*?
- **What is the conditional probability of *prepared*, given *study* and *smart*?**

$$\begin{aligned} p(\text{prepared} | \text{smart}, \text{study}) &= p(\text{prepared}, \text{smart}, \text{study}) / p(\text{smart}, \text{study}) \\ &= .432 / (.432 + .048) \\ &= \mathbf{0.9} \end{aligned}$$

Independence



- When variables don't affect each others' probabilities, they are **independent**; we can easily compute their joint & conditional probability:

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

- {moonPhase, lightLevel} *might* be independent of {burglary, alarm, earthquake}
 - Maybe not: burglars may be more active during a new moon because darkness hides their activity
 - But if we know light level, moon phase doesn't affect whether we are burglarized
 - If burglarized, light level doesn't affect if alarm goes off
- Need a more complex notion of independence and methods for reasoning about the relationships



Exercise: Independence

$p(\text{smart} \wedge \text{study} \wedge \text{prepared})$	smart		\neg smart	
	study	\neg study	study	\neg study
prepared	.432	.16	.084	.008
\neg prepared	.048	.16	.036	.072

Queries:

- Q1: Is *smart* independent of *study*?
- Q2: Is *prepared* independent of *study*?

How can we tell?



Exercise: Independence

$p(\text{smart} \wedge \text{study} \wedge \text{prepared})$	smart		$\neg\text{smart}$	
	study	$\neg\text{study}$	study	$\neg\text{study}$
prepared	.432	.16	.084	.008
$\neg\text{prepared}$.048	.16	.036	.072

Q1: Is *smart* independent of *study*?

- You might have some intuitive beliefs based on your experience
- You can also check the data

Which way to answer this is better?



Exercise: Independence

$p(\text{smart} \wedge \text{study} \wedge \text{prepared})$	smart		\neg smart	
	study	\neg study	study	\neg study
prepared	.432	.16	.084	.008
\neg prepared	.048	.16	.036	.072

Q1: Is *smart* independent of *study*?

Q1 true iff $p(\text{smart} | \text{study}) == p(\text{smart})$

$$p(\text{smart}) = .432 + 0.048 + .16 + .16 = \mathbf{0.8}$$

$$p(\text{smart} | \text{study}) = p(\text{smart}, \text{study}) / p(\text{study}) \\ = (.432 + .048) / .6 = 0.48 / .6 = \mathbf{0.8}$$

$0.8 == 0.8 \therefore$ smart is independent of study



Exercise: Independence

$p(\text{smart} \wedge \text{study} \wedge \text{prep})$	smart		\neg smart	
	study	\neg study	study	\neg study
prepared	.432	.16	.084	.008
\neg prepared	.048	.16	.036	.072

Q2: Is *prepared* independent of *study*?

- What is prepared?
- Q2 true iff



Exercise: Independence

$p(\text{smart} \wedge \text{study} \wedge \text{prep})$	smart		\neg smart	
	study	\neg study	study	\neg study
prepared	.432	.16	.084	.008
\neg prepared	.048	.16	.036	.072

Q2: Is *prepared* independent of *study*?

Q2 true iff $p(\text{prepared} | \text{study}) == p(\text{prepared})$

$$p(\text{prepared}) = .432 + .16 + .84 + .008 = .684$$

$$p(\text{prepared} | \text{study}) = p(\text{prepared}, \text{study}) / p(\text{study})$$

$$= (.432 + .084) / .6 = .86$$

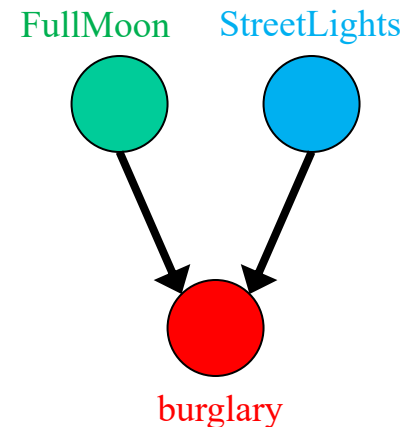
$0.86 \neq 0.684$, \therefore **prepared not independent of study**

Absolute & conditional independence

- Absolute independence:
 - A and B are **independent** if $P(A \wedge B) = P(A) * P(B)$;
equivalently, $P(A) = P(A | B)$ and $P(B) = P(B | A)$
- A and B are **conditionally independent** given C if
 - $P(A \wedge B | C) = P(A | C) * P(B | C)$
- This lets us decompose the joint distribution:
 - $P(A \wedge B \wedge C) = P(A | C) * P(B | C) * P(C)$
- Moon-Phase and Burglary are **conditionally independent given** Light-Level
- Conditional independence is weaker than absolute independence, but useful in decomposing full joint probability distribution

Conditional independence

- Intuitive understanding: conditional independence often comes from **causal relations**
 - FullMoon causally affects LightLevel at night as does StreetLights
- For our burglary scenario, FullMoon doesn't affect anything else
- Knowing *LightLevel*, we can ignore *FullMoon* and *StreetLights* when predicting if alarm suggests **Burglary**



Bayes' rule

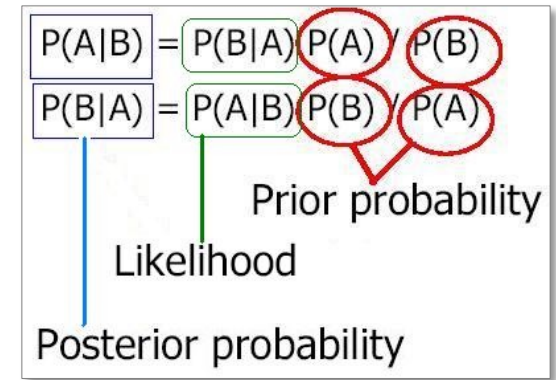
Derived from the product rule:

- $P(A, B) = P(A | B) * P(B)$ *# from definition of conditional probability*
- $P(B, A) = P(B | A) * P(A)$ *# from definition of conditional probability*
- $P(A, B) = P(B, A)$ *# since order is not important*

So...

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

relates $P(A|B)$
and $P(B|A)$



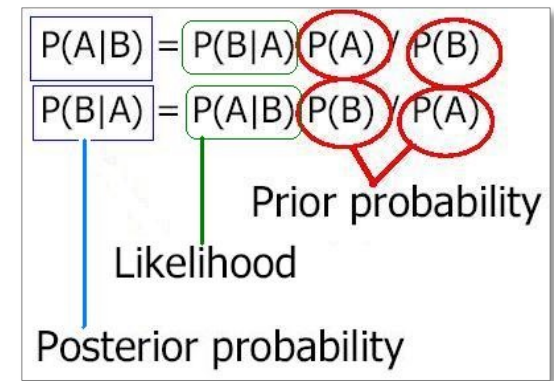
Useful for diagnosis!

- *C is a cause, E is an effect:*

- $P(C|E) = P(E|C) * P(C) / P(E)$

- **Useful for diagnosis:**

- E are (observed) effects and C are (hidden) causes,
 - Often have model for how causes lead to effects $P(E|C)$
 - May also have info (based on experience) on frequency of causes ($P(C)$)
 - Which allows us to reason abductively from effects to causes ($P(C|E)$)
 - Recall, abductive reasoning: from $A \Rightarrow B$ and B , infer (maybe?) A



Ex: meningitis and stiff neck

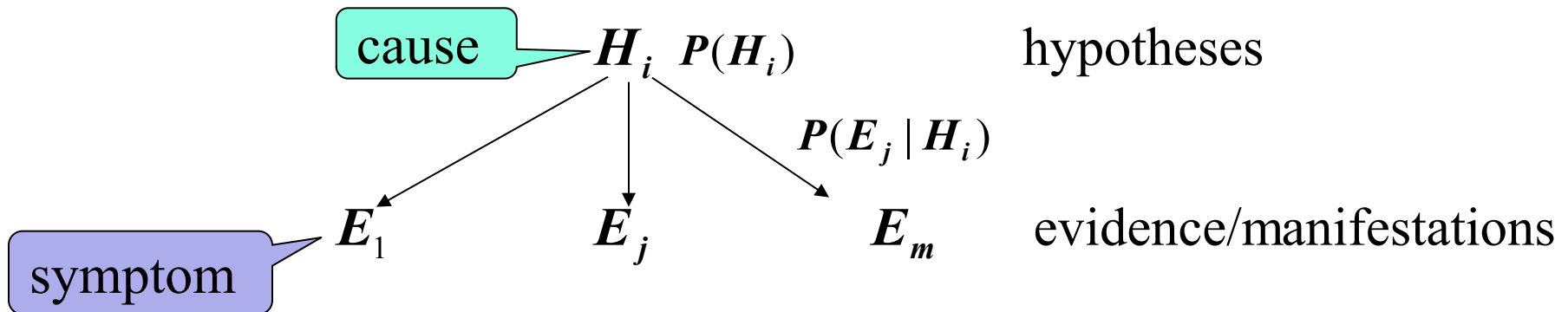
cause

symptom

- **Meningitis (M)** can cause **stiff neck (S)**, though there are other causes too
- Use *S* as a *diagnostic symptom* and estimate $p(M|S)$
- Studies can estimate $p(M)$, $p(S)$ & $p(S|M)$, e.g. $p(S|M)=0.7$, $p(S)=0.01$, $p(M)=0.00002$
- Harder to directly gather data on $p(M|S)$
- Applying Bayes' Rule:
$$p(M|S) = p(S|M) * p(M) / p(S) = 0.0014$$

From multiple evidence to a cause

- In the setting of diagnostic/evidential reasoning



- Know prior probability of hypothesis $P(H_i)$
- conditional probability $P(E_j | H_i)$
- Want to compute the *posterior probability* $P(H_i | E_j)$

- Bayes' s theorem:

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

Bayesian diagnostic reasoning

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

Bayesian diagnostic reasoning (2)

- Prior vs. posterior probability
 - Prior: probability before we know the evidence, e.g., 0.005 for having COVID)
 - Posterior: probability after knowing evidence, e.g., 0.9 if patient tests positive for COVID
- Goal: find hypothesis H_i with highest posterior
 - $\text{Max}_i P(H_i | E_1, \dots, E_l)$
- Requires knowing joint evidence probabilities
$$P(H_i | E_1 \dots E_m) = P(E_1 \dots E_m | H_i) P(H_i) / P(E_1 \dots E_m)$$
- Having many E_i is a big data collection problem!

Simplifying Bayesian diagnostic reasoning

- Having many E_i is a big data collection problem!
- Two ways to address this
- #1 use conditional independence to effect “causal reasoning” and eliminate some E_i
 - Knowing *LightLevel*, we can ignore *FullMoon* and *StreetLights* when predicting if alarm suggests *Burglary*
 - More on this later
- #2 Use a Naïve Bayes approximation that assumes evidence variables are all mutually independent

Naïve Bayesian diagnostic reasoning

- Bayes' rule:

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

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

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

- Easy to compute since we ignore evidence dependence
- Over-simplification for many reasons, but often used as a simple baseline



Summary

- Probability a rigorous formalism for uncertain knowledge
- **Joint probability distribution** specifies probability of every **atomic event**
- Answer queries by summing over atomic events
- Must reduce joint size for non-trivial domains
- **Bayes rule**: compute from known conditional probabilities, usually in causal direction
- **Independence & conditional independence** provide tools
- Next: Bayesian belief networks