

# **Probability Theory for Inference**

# Discrete random variables

---

- A **random variable** can take on one of a set of different values, each with an associated probability. Its value at a particular time is **subject to random variation**.
  - **Discrete** random variables take on one of a discrete (often finite) range of values
  - Domain values must be **exhaustive** and **mutually exclusive**
- For us, random variables will have a discrete, countable (usually finite) domain of **arbitrary values**.
  - Mathematical statistics usually calls these **random elements**
  - **Example: Weather** is a discrete random variable with domain {sunny, rain, cloudy, snow}.
  - **Example: A Boolean random variable** has the domain {true,false},

## A word on notation

---

Assume *Weather* is a discrete random variable with domain {sunny, rain, cloudy, snow}.

- |                                |             |                      |
|--------------------------------|-------------|----------------------|
| • <i>Weather = sunny</i>       | abbreviated | <i>sunny</i>         |
| • <i>P(Weather=sunny)=0.72</i> | abbreviated | <i>P(sunny)=0.72</i> |
| • <i>Cavity = true</i>         | abbreviated | <i>cavity</i>        |
| • <i>Cavity = false</i>        | abbreviated | $\neg$ <i>cavity</i> |

Vector notation:

- Fix order of domain elements:  
*<sunny,rain,cloudy,snow>*
- Specify the probability mass function (pmf) by a vector:  
*P(Weather) = <0.72,0.1,0.08,0.1>*

### 13.2.3 Probability Axioms

- The axiomatization of probability theory by Kolmogorov (1933) based on three simple axioms
1. For any proposition  $a$  the probability is in between 0 and 1:  
$$0 \leq P(a) \leq 1$$
  2. Necessarily true (i.e., valid) propositions have probability 1 and necessarily false (i.e., unsatisfiable) propositions have probability 0:  
$$P(\text{true}) = 1 \quad P(\text{false}) = 0$$
  3. The probability of a disjunction is given by the *inclusion-exclusion principle*  
$$P(a \vee b) = P(a) + P(b) - P(a \wedge b)$$



# Probability Theory

- **Random variables**
  - Domain
- **Atomic event**: complete specification of state
- **Prior probability**: degree of belief without any other evidence
- **Joint probability**: matrix of combined probabilities of a set of variables
- Alarm, Burglary, Earthquake
  - Boolean (like these), discrete, continuous
- $\text{Alarm}=\text{True} \wedge \text{Burglary}=\text{True} \wedge \text{Earthquake}=\text{False}$   
 $\text{alarm} \wedge \text{burglary} \wedge \neg \text{earthquake}$
- $P(\text{Burglary}) = .1$
- $P(\text{Alarm}, \text{Burglary}) =$

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



# Probability Theory: Definitions

- **Computing conditional prob:**

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

## Bayes' Rule & Diagnosis

$$\underset{\text{Posterior}}{P(a|b)} = \frac{\overset{\text{Likelihood}}{P(b|a)} * \overset{\text{Prior}}{P(a)}}{\underset{\text{Normalization}}{P(b)}}$$

- Useful for assessing **diagnostic** probability from **causal** probability:

$$P(\text{Cause}|\text{Effect}) = \frac{P(\text{Effect}|\text{Cause}) * P(\text{Cause})}{P(\text{Effect})}$$

# Probability Summary

Conditional probability

$$P(x|y) = \frac{P(x, y)}{P(y)}$$

Product rule

$$P(x, y) = P(x|y)P(y)$$

Chain rule

$$\begin{aligned} P(X_1, X_2, \dots, X_n) &= P(X_1)P(X_2|X_1)P(X_3|X_1, X_2) \dots \\ &= \prod_{i=1}^n P(X_i|X_1, \dots, X_{i-1}) \end{aligned}$$

X and Y independent if and only if:  $\forall x, y : P(x, y) = P(x)P(y)$

X and Y are conditionally independent given Z if and only if:

$$\forall x, y, z : P(x, y|z) = P(x|z)P(y|z)$$

$$X \perp\!\!\!\perp Y|Z$$

# Try It...

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

- $P(\text{alarm} \mid \text{burglary}) = ??$
- $P(\text{burglary} \mid \text{alarm}) = ??$
- $P(\text{burglary} \wedge \text{alarm}) = ??$
- $P(\text{alarm}) = ??$

- **Computing conditional prob:**

- $P(a \mid b) = P(a \wedge b) / P(b)$
- $P(b)$ : **normalizing** constant

- **Product rule:**

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

- **Marginalizing:**

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



# Probability Theory (cont.)

- **Conditional probability:**  
probability 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**)
- $P(\text{burglary} | \text{alarm}) = .47$   
 $P(\text{alarm} | \text{burglary}) = .9$
- $P(\text{burglary} | \text{alarm}) =$   
 $P(\text{burglary} \wedge \text{alarm}) / P(\text{alarm})$   
 $= .09 / .19 = .47$
- $P(\text{burglary} \wedge \text{alarm}) =$   
 $P(\text{burglary} | \text{alarm}) P(\text{alarm}) =$   
 $.47 * .19 = .09$
- $P(\text{alarm}) =$   
 $P(\text{alarm} \wedge \text{burglary}) +$   
 $P(\text{alarm} \wedge \neg \text{burglary}) =$   
 $.09 + .1 = .19$

## Bayes Theorem

### Application

Guilty or not?

*A person is put in front of a jury. The jury finds the defendant guilty in 98% of the cases in which the defendant has committed a crime, and it finds the defendant not guilty in only 91% of the cases in which the defendant has not committed a crime. Furthermore, only .008 of the entire population has committed a crime.*

If a random person is found guilty by the jury, what's more likely: criminal or not?

## Bayes Theorem Application

### Guilty or not?

*A person is put in front of a jury. The jury finds the defendant guilty in 98% of the cases in which the defendant has committed a crime, and it finds the defendant not guilty in only 97% of the cases in which the defendant has not committed a crime.*

*Furthermore, only .008 of the entire population has committed a crime.*

$$P(\text{criminal}) = 0.008$$

$$P(\neg \text{criminal}) = 0.992$$

$$P(\text{guilty}|\text{criminal}) = \underline{0.98}$$

$$P(\neg \text{guilty}|\text{criminal}) = \underline{0.02}$$

$$P(\text{guilty}|\neg \text{criminal}) = 0.03$$

$$P(\neg \text{guilty}|\neg \text{criminal}) = \underline{0.97}$$

If a random person is found guilty by the jury, what's more likely: criminal or not?

which is bigger?  $P(\text{criminal}|\text{guilty})$  or  $P(\neg \text{criminal}|\text{guilty})$ ?

## Probabilities

### Bayes Rule

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

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

$$P(b|a)P(a) = P(a|b)P(b)$$

$$P(\underline{b}|\underline{a}) = \frac{P(a|b)P(b)}{P(a)}$$



## Bayes Theorem Application

Guilty or not?

$$P(\text{criminal}) = 0.008$$

$$P(\neg \text{criminal}) = 0.992$$

$$P(\text{guilty}|\text{criminal}) = 0.98$$

$$P(\neg \text{guilty}|\text{criminal}) = 0.03$$

$$P(\text{guilty}|\neg \text{criminal}) = 0.02$$

$$P(\neg \text{guilty}|\neg \text{criminal}) = 0.97$$

If a random person is found guilty by the jury, what's more likely: criminal or not?

which is bigger?  $P(\text{criminal}|\text{guilty})$  or  $P(\neg \text{criminal}|\text{guilty})$ ?

$$\text{— } \underline{P(\text{criminal}|\text{guilty})} = \frac{P(\text{guilty}|\text{criminal})P(\text{criminal})}{P(\text{guilty})}$$

$$\text{— } P(\neg \text{criminal}|\text{guilty}) = \frac{P(\text{guilty}|\neg \text{criminal})P(\neg \text{criminal})}{P(\text{guilty})}$$

# Calculating Conditional Probabilities

College students were asked if they have ever cheated on an exam. Results were broken down by gender.

		Cheated on College Exam?		
		Yes	No	Total
Gender	Male	.32	.22	.54
	Female	.28	.18	.46
	Total	.60	.40	1.00

● Question: Given that a person has cheated, what is the probability he is male?

● Answer: 
$$P(\text{Male}|\text{Cheater}) = \frac{P(\text{Male} \cap \text{Cheater})}{P(\text{Cheater})}$$
$$= \frac{.32}{.60} = .5333$$



	Right-handed	Left-handed	Total
Male	0.41	0.08	0.49
Female	0.45	0.06	0.51
Total	0.86	0.14	1

Find the probability that a randomly selected person is:

- (a) a male given that she is right-handed;
- (b) right-handed given that she is a male;
- (c) a female given that she is left-handed.
- (d) Are the events *being a female* and *being left-handed* independent? Justify.

$$\underline{a)} \quad P(M | R) = \frac{P(M \cap R)}{P(R)} = \frac{0.41}{0.86} \approx 0.477$$

$$\underline{b)} \quad P(R | M) = \frac{P(R \cap M)}{P(M)} = \frac{0.41}{0.49} \approx 0.837$$

$$\underline{b)} \quad P(R|M) = \frac{P(R \cap M)}{P(M)} = \frac{0.41}{0.49} \approx 0.837$$

$$\underline{c)} \quad P(F|L) = \frac{P(F \cap L)}{P(L)} = \frac{0.06}{0.14} \approx 0.429$$

$$\underline{d)} \quad \left. \begin{array}{l} P(F|L) \approx 0.429 \\ P(F) = 0.51 \end{array} \right\} \neq \Rightarrow F \text{ and } L \text{ are } \underline{\text{not}} \text{ independent}$$

# Joint probability distribution

- Probability assignment to all combinations of values of random variables (i.e. all elementary events)

	toothache	$\neg$ toothache
cavity	0.04	0.06
$\neg$ cavity	0.01	0.89

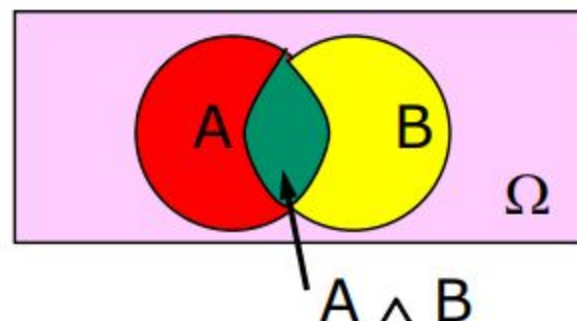


- The sum of the entries in this table has to be 1
- *Every question about a domain can be answered by the joint distribution*
- Probability of a proposition is the sum of the probabilities of elementary events in which it holds
  - $P(\text{cavity}) = 0.1$  [marginal of row 1]
  - $P(\text{toothache}) = 0.05$  [marginal of toothache column]



# Conditional Probability

	toothache	$\neg$ toothache
cavity	0.04	0.06
$\neg$ cavity	0.01	0.89



- $P(\text{cavity})=0.1$  and  $P(\text{cavity} \wedge \text{toothache})=0.04$  are both *prior* (unconditional) probabilities
- Once the agent has new evidence concerning a *previously unknown* random variable, e.g. Toothache, we can specify a *posterior* (conditional) probability e.g.  $P(\text{cavity} \mid \text{Toothache}=\text{true})$

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

[Probability of  $a$  with the Universe  $\Omega$  restricted to  $b$ ]

→ The new information restricts the set of possible worlds  $\omega_i$  consistent with it, so **changes the probability**.

- So  $P(\text{cavity} \mid \text{toothache}) = 0.04 / 0.05 = 0.8$



# Conditional Probability (continued)

---

- **Definition of Conditional Probability:**

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

- **Product rule gives an alternative formulation:**

$$\begin{aligned} P(a \wedge b) &= P(a \mid b) * P(b) \\ &= P(b \mid a) * P(a) \end{aligned}$$

- **A general version holds for whole distributions:**

$$P(\textit{Weather}, \textit{Cavity}) = P(\textit{Weather} \mid \textit{Cavity}) * P(\textit{Cavity})$$

- **Chain rule** is derived by successive application of product rule:

$$\begin{aligned} P(A, B, C, D, E) &= P(A \mid B, C, D, E) P(B, C, D, E) \\ &= P(A \mid B, C, D, E) P(B \mid C, D, E) P(C, D, E) \\ &= \dots \\ &= P(A \mid B, C, D, E) P(B \mid C, D, E) P(C \mid D, E) P(D \mid E) P(E) \end{aligned}$$

# Probabilistic Inference

---

- **Probabilistic inference:** the computation
  - from *observed evidence*
  - of *posterior probabilities*
  - for *query propositions*.
- We use the **full joint distribution** as the “knowledge base” from which answers to questions may be derived.
- Ex: three Boolean variables *Toothache (T)*, *Cavity (C)*, *ShowsOnXRay (X)*

	t		$\neg t$	
	x	$\neg x$	x	$\neg x$
c	0.108	0.012	0.072	0.008
$\neg c$	0.016	0.064	0.144	0.576

- Probabilities in joint distribution sum to 1



# Probabilistic Inference II

---

	t		$\neg t$	
	x	$\neg x$	x	$\neg x$
c	0.108	0.012	0.072	0.008
$\neg c$	0.016	0.064	0.144	0.576

- Probability of any proposition computed by finding atomic events where proposition is true and adding their probabilities
  - $P(\text{cavity} \vee \text{toothache})$   
 $= 0.108 + 0.012 + 0.072 + 0.008 + 0.016 + 0.064$   
 $= 0.28$
  - $P(\text{cavity})$   
 $= 0.108 + 0.012 + 0.072 + 0.008$   
 $= 0.2$
- $P(\text{cavity})$  is called a marginal probability and the process of computing this is called marginalization

# Probabilistic Inference III

---

	t		$\neg t$	
	x	$\neg x$	x	$\neg x$
c	0.108	0.012	0.072	0.008
$\neg c$	0.016	0.064	0.144	0.576

- Can also compute conditional probabilities.
- $P(\neg \text{cavity} \mid \text{toothache})$   
     $= P(\neg \text{cavity} \wedge \text{toothache}) / P(\text{toothache})$   
     $= (0.016 + 0.064) / (0.108 + 0.012 + 0.016 + 0.064)$   
     $= 0.4$
- Denominator is viewed as a *normalization constant*:
  - Stays constant no matter what the value of Cavity is.  
(Book uses  $\alpha$  to denote normalization constant  $1/P(X)$ , for random variable  $X$ .)


## 13.3 Inference Using Full Joint Distribution



	toothache		¬toothache	
	catch	¬catch	catch	¬catch
cavity	0.108	0.012	0.072	0.008
¬cavity	0.016	0.064	0.144	0.576

- E.g., there are six atomic events for  $cavity \vee toothache$ :  
 $0.108 + 0.012 + 0.072 + 0.008 + 0.016 + 0.064 = 0.28$
- Extracting the distribution over a variable (or some subset of variables), *marginal probability*, is attained by adding the entries in the corresponding rows or columns
- E.g.,  $P(cavity) = 0.108 + 0.012 + 0.072 + 0.008 = 0.2$
- We can write the following general marginalization (summing out) rule for any sets of variables  $\mathbf{Y}$  and  $\mathbf{Z}$ :

$$P(\mathbf{Y}) = \sum_{\mathbf{z} \in \mathbf{Z}} P(\mathbf{Y}, \mathbf{z})$$



	toothache		¬toothache	
	catch	¬catch	catch	¬catch
cavity	0.108	0.012	0.072	0.008
¬cavity	0.016	0.064	0.144	0.576

- Computing a conditional probability

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

$$P(\text{cavity} \wedge \text{toothache}) / P(\text{toothache}) =$$

$$(0.108 + 0.012) / (0.108 + 0.012 + 0.016 + 0.064) =$$

$$0.12 / 0.2 = 0.6$$

- Respectively

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

$$(0.016 + 0.064) / 0.2 = 0.4$$

- The two probabilities sum up to one, as they should



## 13.4 Independence

- If we expand the previous example with a fourth random variable *Weather*, which has four possible values, we have to copy the table of joint probabilities four times to have 32 entries together

- Dental problems have no influence on the weather, hence:

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

- By this observation and product rule

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

# 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 still useful in decomposing the full joint probability distribution



# Bayes' Rule

- Two ways to factor a joint distribution over two variables:

$$P(x, y) = P(x|y)P(y) = P(y|x)P(x)$$

That's my rule!

- Dividing, we get:

$$P(x|y) = \frac{P(y|x)}{P(y)}P(x)$$



- Why is this at all helpful?
  - Lets us build a conditional from its reverse
  - Often one conditional is tricky but the other one is simple
  - Foundation of many systems we'll see later
- In the running for most important AI equation!

# Bayes' Rule & Diagnosis

$$\underset{\text{Posterior}}{P(a|b)} = \frac{\overset{\text{Likelihood}}{P(b|a)} * \overset{\text{Prior}}{P(a)}}{\underset{\text{Normalization}}{P(b)}}$$

Useful for assessing **diagnostic** probability from **causal** probability:

$$P(\text{Cause}|\text{Effect}) = \frac{P(\text{Effect}|\text{Cause}) * P(\text{Cause})}{P(\text{Effect})}$$

# Bayes' Rule For Diagnosis II

---

$$P(\text{Disease} \mid \text{Symptom}) = \frac{P(\text{Symptom} \mid \text{Disease}) * P(\text{Disease})}{P(\text{Symptom})}$$

Imagine:

- disease = TB, symptom = coughing
- $P(\text{disease} \mid \text{symptom})$  is different in TB-indicated country vs. USA
- $P(\text{symptom} \mid \text{disease})$  should be the same
  - It is more widely useful to learn  $P(\text{symptom} \mid \text{disease})$
- What about  $P(\text{symptom})$ ?
  - Use *conditioning* (next slide)
  - For determining, e.g., the *most likely* disease given the symptom, we can just ignore  $P(\text{symptom})$ !!! (see slide 35)



# Conditioning

- Idea:** Use *conditional probabilities* instead of joint probabilities

$$\begin{aligned} P(a) &= P(a \wedge b) + P(a \wedge \neg b) \\ &= P(a \mid b) * P(b) + P(a \mid \neg b) * P(\neg b) \end{aligned}$$

Here:

$$\begin{aligned} P(\text{symptom}) &= P(\text{symptom} \mid \text{disease}) * P(\text{disease}) \\ &\quad P(\text{symptom} \mid \neg \text{disease}) * P(\neg \text{disease}) \end{aligned}$$

- More generally:  $P(Y) = \sum_z P(Y|z) * P(z)$
- Marginalization and conditioning are useful rules for derivations involving probability expressions.

# Conditional Independence

BUT *absolute* independence is rare

Dentistry is a large field with hundreds of variables, one of which are independent. What to do?

A and B are conditionally independent given C iff

- $P(A | B, C) = P(A | C)$
- $P(B | A, C) = P(B | C)$
- $P(A \wedge B | C) = P(A | C) * P(B | C)$

Toothache (T), Spot in Xray (X), Cavity (C)

- None of these are independent of the other two
- But ***T and X are conditionally independent given C***





# Conditional Independence

---

- $P(\text{Toothache}, \text{Cavity}, \text{Catch})$
- If I have a cavity, the probability that the probe catches in it doesn't depend on whether I have a toothache:
  - $P(+\text{catch} \mid +\text{toothache}, +\text{cavity}) = P(+\text{catch} \mid +\text{cavity})$
- The same independence holds if I don't have a cavity:
  - $P(+\text{catch} \mid +\text{toothache}, -\text{cavity}) = P(+\text{catch} \mid -\text{cavity})$
- Catch is *conditionally independent* of Toothache given Cavity:
  - $P(\text{Catch} \mid \text{Toothache}, \text{Cavity}) = P(\text{Catch} \mid \text{Cavity})$
- Equivalent statements:
  - $P(\text{Toothache} \mid \text{Catch}, \text{Cavity}) = P(\text{Toothache} \mid \text{Cavity})$
  - $P(\text{Toothache}, \text{Catch} \mid \text{Cavity}) = P(\text{Toothache} \mid \text{Cavity}) P(\text{Catch} \mid \text{Cavity})$
  - One can be derived from the other easily



## Conditional Independence II *WHY??*

If I have a cavity, the probability that the XRay shows a spot doesn't depend on whether I have a toothache (and vice versa)

$$P(X|T,C) = P(X|C)$$

From which follows:

$$P(T|X,C) = P(T|C) \quad \text{and} \quad P(T,X|C) = P(T|C) * P(X|C)$$

By the chain rule), given conditional independence:

$$\begin{aligned} P(T,X,C) &= P(T|X,C) * P(X,C) = P(T|X,C) * P(X|C) * P(C) \\ &= P(T|C) * P(X|C) * P(C) \end{aligned}$$

$P(\text{Toothache}, \text{Cavity}, \text{Xray})$  has  $2^3 - 1 = 7$  independent entries

Given conditional independence, chain rule yields  
 $2 + 2 + 1 = 5$  independent numbers

# Conditional Independence III

---

- 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.*

# Exercise: Inference from the Joint

$p(\text{smart} \wedge \text{study} \wedge \text{prep})$	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*?
- Save these answers for later! 😊

# Exercise: Independence

$p(\text{smart} \wedge \text{study} \wedge \text{prep})$	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:
  - Is *smart* independent of *study*?
  - Is *prepared* independent of *study*?



# Exercise: Conditional Independence

$p(\text{smart} \wedge \text{study} \wedge \text{prep})$	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:
  - Is *smart* conditionally independent of *prepared*, given *study*?
  - Is *study* conditionally independent of *prepared*, given *smart*?