

Representing Uncertainty for Probabilistic Inference

Chapter 13

1

Uncertainty in the World

- An agent can often be uncertain about the state of the domain since there is often ambiguity and uncertainty
- Plausible/**probabilistic inference**
 - I've got this evidence; what's the chance that this conclusion is true?
 - I've got a sore neck; how likely am I to have meningitis?
 - A mammogram test is positive; what's the probability that the patient has breast cancer?

2

Probabilistic Inference

How do we use probabilities in AI?

- You wake up with a headache
- Do you have the flu?
- H = headache, F = flu



Logical Inference: if H then F
(but the world is usually not this simple)

Statistical Inference: compute the probability of a query/diagnosis/decision given (i.e., conditioned on) evidence/symptom/observation, i.e., $P(F | H)$

[Example from Andrew Moore]

3

Uncertainty

- Say we have a rule:
*if toothache **then** problem is cavity*
- But not all patients have toothaches due to cavities, so we could set up rules like:
*if toothache and \neg gum-disease and \neg filling and ...
then problem = cavity*
- This gets complicated; better method:
*if toothache **then** problem is cavity with 0.8 probability*
or $P(\text{cavity} | \text{toothache}) = 0.8$
the probability of cavity is 0.8 given toothache is observed

4

Uncertainty in the World and our Models

- True uncertainty: *rules are probabilistic in nature*
 - quantum mechanics
 - rolling dice, flipping a coin
- Laziness: *too hard to determine exception-less rules*
 - takes too much work to determine *all* of the relevant factors
 - too hard to use the enormous rules that result
- Theoretical ignorance: *don't know all the rules*
 - problem domain has no complete, consistent theory (e.g., medical diagnosis)
- Practical ignorance: *do know all the rules BUT*
 - haven't collected all relevant information for a particular case

6

Logics

Logics are characterized by what they use as "primitives"

Logic	What Exists in World	Knowledge States
Propositional	facts	true/false/unknown
First-Order	facts, objects, relations	true/false/unknown
Temporal	facts, objects, relations, times	true/false/unknown
Probability Theory	facts	degree of belief 0..1
Fuzzy	degree of truth	degree of belief 0..1

7

Probability Theory

- **Probability theory** serves as a formal means for
 - Representing and reasoning with uncertain knowledge
 - Modeling **degrees of belief** in a proposition (event, conclusion, diagnosis, etc.)
- *Probability is the "language" of uncertainty*
 - A key modeling tool in modern AI

8

Sample Space

- A space of **events** in which we assign probabilities
- Events can be binary, multi-valued, or continuous
- Events are **mutually exclusive**
- Examples
 - Coin flip: {head, tail}
 - Die roll: {1, 2, 3, 4, 5, 6}
 - English words: a dictionary
 - High temperature tomorrow: {-100, ..., 100}

10

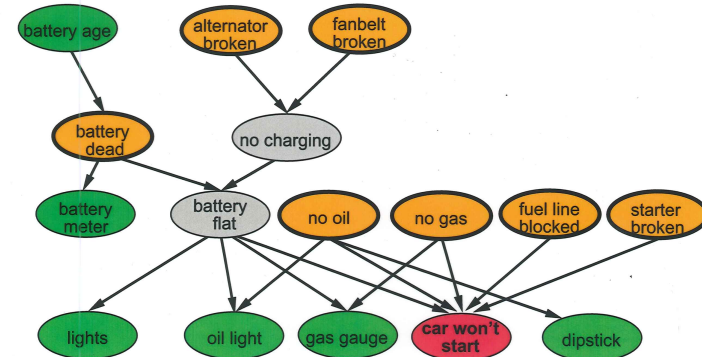
Random Variable

- A variable, X , whose domain is a sample space, and whose value is (somewhat) uncertain
- Examples:
 - X = coin flip outcome
 - X = tomorrow's high temperature
- For a given task, the user defines a set of random variables for describing the world
- **Each variable has a set of mutually exclusive and exhaustive possible values**

11

Example: Car diagnosis

Initial evidence: car won't start
 Testable variables (green), "broken, so fix it" variables (orange)
 Hidden variables (gray) ensure sparse structure, reduce parameters



13

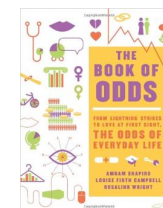
Probability for Discrete Events

- An agent's uncertainty is represented by $P(A=a)$ or simply $P(a)$
 - the agent's degree of belief that variable A takes on value a given no other information related to A
 - a single probability called an **unconditional** or **prior probability**

14

Probability for Discrete Events

- Examples
 - $P(\text{head}) = P(\text{tail}) = 0.5$ fair coin
 - $P(\text{head}) = 0.51, P(\text{tail}) = 0.49$ slightly biased coin
 - $P(\text{first word} = \text{"the"})$ when flipping to a random page in R&N = ?



- Book: *The Book of Odds*

15

Source of Probabilities

- **Frequentists**
 - probabilities come from data
 - if 10 of 100 people tested have a cavity, $P(\text{cavity}) = 0.1$
 - probability means the fraction that would be observed in the limit of infinitely many samples
- **Objectivists**
 - probabilities are real aspects of the world
 - objects have a propensity to behave in certain ways
 - coin has propensity to come up heads with probability 0.5
- **Subjectivists**
 - probabilities characterize an agent's belief
 - have no external physical significance

18

Probability Distributions

Given A is a RV taking values in $\langle a_1, a_2, \dots, a_k \rangle$

e.g., if A is *Sky*, then value is one of $\langle \text{clear}, \text{partly_cloudy}, \text{overcast} \rangle$

- $P(a)$ represents a **single probability** where $A=a$

e.g., if A is *Sky*, then $P(a)$ means any one of $P(\text{clear}), P(\text{partly_cloudy}), P(\text{overcast})$

- $P(A)$ represents a **probability distribution**

- the **set of values**: $\langle P(a_1), P(a_2), \dots, P(a_k) \rangle$

- If A takes n values, then $P(A)$ is a set of n probabilities

e.g., if A is *Sky*, then $P(\text{Sky})$ is the set of probabilities:

$\langle P(\text{clear}), P(\text{partly_cloudy}), P(\text{overcast}) \rangle$

- Property: $\sum P(a_i) = P(a_1) + P(a_2) + \dots + P(a_k) = 1$

- sum over all values in the domain of variable A is 1 because the **domain is mutually exclusive and exhaustive**

19

Probability Table

- *Weather*

sunny	cloudy	rainy
200/365	100/365	65/365

- $P(\text{Weather} = \text{sunny}) = P(\text{sunny}) = 200/365$
- $P(\text{Weather}) = \langle 200/365, 100/365, 65/365 \rangle$
- We'll obtain the probabilities by counting frequencies from data

20

The Axioms of Probability

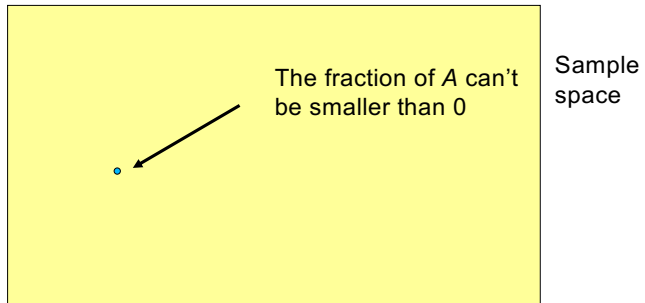
1. $0 \leq P(A) \leq 1$
2. $P(\text{true}) = 1, P(\text{false}) = 0$
3. $P(A \vee B) = P(A) + P(B) - P(A \wedge B)$

Note: Here $P(A)$ means $P(A=a)$ for some value a and $P(A \vee B)$ means $P(A=a \vee B=b)$

21

The Axioms of Probability

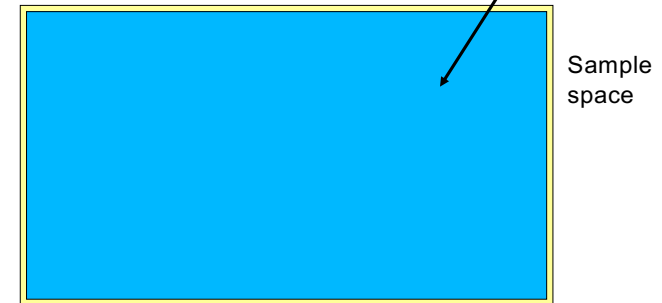
- $0 \leq P(A) \leq 1$
- $P(\text{true}) = 1, P(\text{false}) = 0$
- $P(A \vee B) = P(A) + P(B) - P(A \wedge B)$



22

The Axioms of Probability

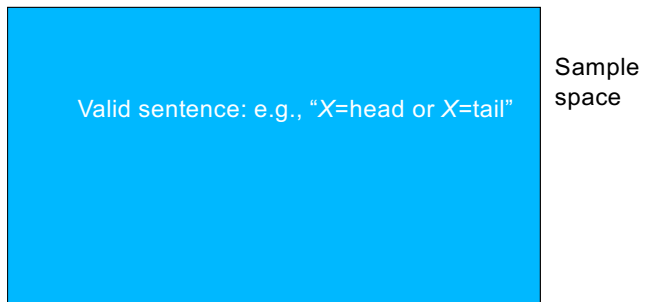
- $0 \leq P(A) \leq 1$
- $P(\text{true}) = 1, P(\text{false}) = 0$
- $P(A \vee B) = P(A) + P(B) - P(A \wedge B)$



23

The Axioms of Probability

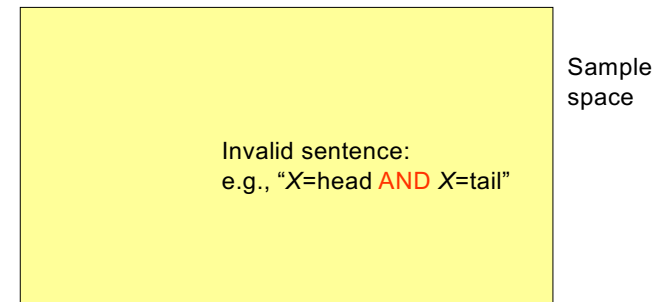
- $0 \leq P(A) \leq 1$
- $P(\text{true}) = 1, P(\text{false}) = 0$
- $P(A \vee B) = P(A) + P(B) - P(A \wedge B)$



24

The Axioms of Probability

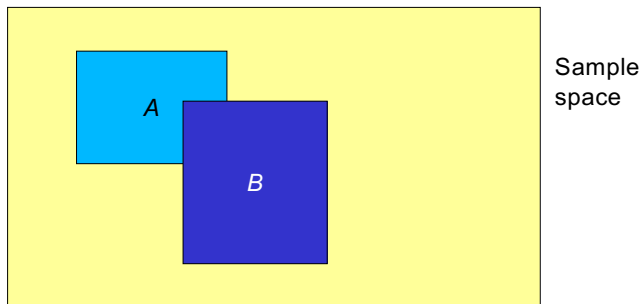
- $0 \leq P(A) \leq 1$
- $P(\text{true}) = 1, P(\text{false}) = 0$
- $P(A \vee B) = P(A) + P(B) - P(A \wedge B)$



25

The Axioms of Probability

- $0 \leq P(A) \leq 1$
- $P(\text{true}) = 1, P(\text{false}) = 0$
- $P(A \vee B) = P(A) + P(B) - P(A \wedge B)$



26

Some Theorems Derived from the Axioms

- $P(\neg A) = 1 - P(A)$
- If A can take k different values a_1, \dots, a_k :

$$P(A=a_1) + \dots + P(A=a_k) = 1$$
- $P(B) = P(B \wedge \neg A) + P(B \wedge A)$, if A is a binary event
- $P(B) = \sum_{i=1 \dots k} P(B \wedge A=a_i)$, if A can take k values

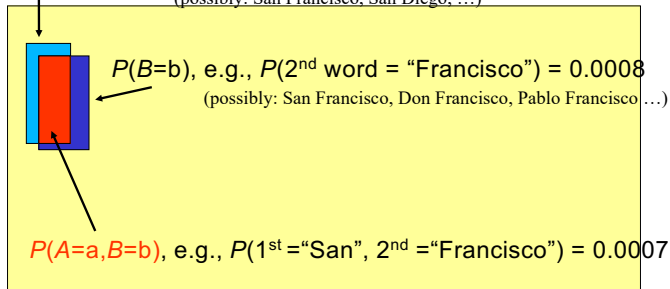
Called **Addition** or **Conditioning rule**

27

Joint Probability

- The **joint probability** $P(A=a, B=b)$ is shorthand for $P(A=a \wedge B=b)$, i.e., the probability of *both* $A=a$ and $B=b$ happening

$P(A=a)$, e.g., $P(1^{\text{st}} \text{ word on a random page} = \text{"San"}) = 0.001$
 (possibly: San Francisco, San Diego, ...)



28

Full Joint Probability Distribution (FJPD)

		Weather		
		<i>sunny</i>	<i>cloudy</i>	<i>rainy</i>
Temp	<i>hot</i>	150/365	40/365	5/365
	<i>cold</i>	50/365	60/365	60/365

- $P(\text{Temp}=\text{hot}, \text{Weather}=\text{rainy}) = P(\text{hot}, \text{rainy}) = 5/365 = 0.014$
- The **full joint probability distribution** table for n random variables, each taking k values, has k^n entries

29

Full Joint Probability Distribution (FJPD)

<i>Bird</i>	<i>Flier</i>	<i>Young</i>	Probability
T	T	T	0.0
T	T	F	0.2
T	F	T	0.04
T	F	F	0.01
F	T	T	0.01
F	T	F	0.01
F	F	T	0.23
F	F	F	0.5

3 Boolean random variables $\Rightarrow 2^3 - 1 = 7$ “degrees of freedom” (DOF) or “independent values”

Sums to 1

30

Computing from the FJPD

- **Marginal Probabilities**

- $P(\text{Bird}=\text{T}) = P(\text{bird}) = 0.0 + 0.2 + 0.04 + 0.01 = 0.25$

- $P(\text{bird}, \neg \text{flier}) = 0.04 + 0.01 = 0.05$

- $P(\text{bird} \vee \text{flier}) = 0.0 + 0.2 + 0.04 + 0.01 + 0.01 + 0.01 = 0.27$

- Sum over all other variables

- **“Summing Out”**

- **“Marginalization”**

31

Unconditional / Prior Probability

- One’s uncertainty or original assumption about an event *prior* to having any data about it **or anything else** in the domain
- $P(\text{Coin} = \text{heads}) = 0.5$
- $P(\text{Bird} = \text{T}) = 0.0 + 0.2 + 0.04 + 0.01 = 0.22$
- Compute from the FJPD by marginalization

32

Marginal Probability

		<i>Weather</i>		
<i>Temp</i>		<i>sunny</i>	<i>cloudy</i>	<i>rainy</i>
	<i>hot</i>	150/365	40/365	5/365
	<i>cold</i>	50/365	60/365	60/365
Σ		200/365	100/365	65/365

$P(\text{Weather}) = \langle 200/365, 100/365, 65/365 \rangle$

Probability **distribution** for r.v. *Weather*

The name comes from the old days when the sums were written in the margin of a page

33

Marginal Probability

		Weather			Σ
		sunny	cloudy	rainy	
Temp	hot	150/365	40/365	5/365	195/365
	cold	50/365	60/365	60/365	170/365

$$P(\text{Temp}) = \langle 195/365, 170/365 \rangle$$

This is nothing but $P(B) = \sum_{i=1 \dots k} P(B \wedge A=a_i)$,
where A can take k values

34

Conditional Probability

Conditional probabilities

- formalizes the process of accumulating evidence and updating probabilities based on new evidence
- specifies the belief in a proposition (event, conclusion, diagnosis, etc.) that is *conditioned on* a proposition (evidence, feature, symptom, etc.) being true

- $P(a \mid e)$: **conditional probability** of $A=a$ given $E=e$ evidence is *all that is known true*

$$P(a \mid e) = P(a \wedge e) / P(e) = P(a, e) / P(e)$$

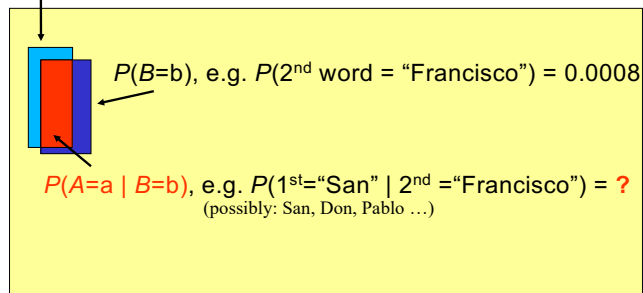
- conditional probability can be viewed as the joint probability $P(a, e)$ normalized by the prior probability, $P(e)$

35

Conditional Probability

The **conditional** probability $P(A=a \mid B=b)$ is the fraction of time $A=a$, **within the region where** $B=b$

$P(A=a)$, e.g. $P(1^{\text{st}} \text{ word on a random page} = \text{"San"}) = 0.001$

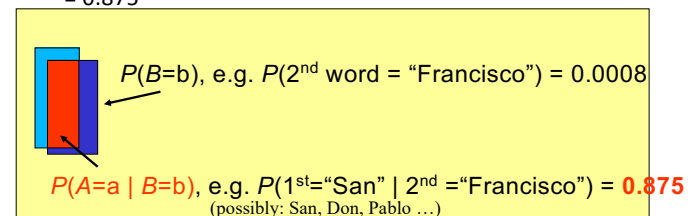


36

Conditional Probability

- $P(\text{san} \mid \text{francisco})$
 $= \#(1^{\text{st}} = \text{s and } 2^{\text{nd}} = \text{f}) / \#(2^{\text{nd}} = \text{f})$
 $= P(\text{san} \wedge \text{francisco}) / P(\text{francisco})$
 $= 0.0007 / 0.0008$
 $= 0.875$

$$\begin{aligned} P(s) &= 0.001 \\ P(f) &= 0.0008 \\ P(s, f) &= 0.0007 \end{aligned}$$



Although "San" is rare and "Francisco" is rare,
given "Francisco" then "San" is quite likely!

37

Conditional Probability

Conditional probabilities behave exactly like standard probabilities; for example:

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

conditional probabilities are between 0 and 1 inclusive

$$P(a_1 \mid e) + P(a_2 \mid e) + \dots + P(a_k \mid e) = 1$$

conditional probabilities sum to 1 where a_1, \dots, a_k are all values in the domain of random variable A

$$P(\neg a \mid e) = 1 - P(a \mid e)$$

negation for Boolean random variable A

38

Computing Conditional Probability

$$P(\neg B \mid F) = ?$$

$$P(F) = ?$$

Note: $P(\neg B \mid F)$ means $P(B=\text{false} \mid F=\text{true})$
and $P(F)$ means $P(F=\text{true})$

40

Full Joint Probability Distribution

Bird (B)	Flier (F)	Young (Y)	Probability
T	T	T	0.0
T	T	F	0.2
T	F	T	0.04
T	F	F	0.01
F	T	T	0.01
F	T	F	0.01
F	F	T	0.23
F	F	F	0.5

3 Boolean random variables $\Rightarrow 2^3 - 1 = 7$
“degrees of freedom” or “independent values”

Sums to 1

41

Computing Conditional Probability

$$\begin{aligned} P(\neg B \mid F) &= P(\neg B, F) / P(F) \\ &= (P(\neg B, F, Y) + P(\neg B, F, \neg Y)) / P(F) \\ &= (0.01 + 0.01) / P(F) \end{aligned}$$

$$\begin{aligned} P(F) &= P(F, B, Y) + P(F, B, \neg Y) + P(F, \neg B, Y) + \\ &\quad P(F, \neg B, \neg Y) \\ &= 0.0 + 0.2 + 0.01 + 0.01 \\ &= 0.22 \end{aligned}$$

Marginalization

42

Computing Conditional Probability

- Instead of using Marginalization to compute $P(F)$, can alternatively use **Normalization**:
- $P(\neg B | F) = .02/P(F)$ from previous slide
- $P(\neg B | F) + P(B | F) = 1$ by definition
- $P(B | F) = P(B, F)/P(F) = (0.0 + 0.2)/P(F)$
- So, $0.02/P(F) + 0.2/P(F) = 1$
- Hence, $P(F) = 0.22$

43

Normalization

- In general, $P(A | B) = \alpha P(A, B)$
where $\alpha = 1/P(B) = 1/(P(A, B) + P(\neg A, B))$
- $P(Q | E_1, \dots, E_k) = \alpha P(Q, E_1, \dots, E_k)$
 $= \alpha \sum_Y P(Q, E_1, \dots, E_k, Y)$

Addition rule

44

Conditional Probability with Multiple Evidence

$$\begin{aligned}
 P(\neg B | F, \neg Y) &= P(\neg B, F, \neg Y) / P(F, \neg Y) \\
 &= P(\neg B, F, \neg Y) / (P(\neg B, F, \neg Y) + P(B, F, \neg Y)) \\
 &= .01 / (.01 + .2) \\
 &= 0.048
 \end{aligned}$$

45

Conditional Probability

- $P(X_1=x_1, \dots, X_k=x_k | X_{k+1}=x_{k+1}, \dots, X_n=x_n) =$
sum of all entries in FJPD where $X_1=x_1, \dots, X_n=x_n$ divided by sum of all entries where $X_{k+1}=x_{k+1}, \dots, X_n=x_n$
- But this means in general we need the *entire* FJPD table, requiring an *exponential number of values* to do probabilistic inference (i.e., compute conditional probabilities)

46

The Chain Rule

- From the definition of conditional probability we have

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

- It also works the other way around:

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

- It works with more than 2 events too:

$$P(A_1, A_2, \dots, A_n) =$$

$$P(A_1) * P(A_2 | A_1) * P(A_3 | A_1, A_2) * \dots$$

$$* P(A_n | A_1, A_2, \dots, A_{n-1})$$

Called
"Product
Rule"

Called "Chain Rule"

48

Probabilistic Reasoning

How do we use probabilities in AI?

- You wake up with a headache
- Do you have the flu?
- H = headache, F = flu



Logical Inference: if H then F
(but the world is usually not this simple)

Statistical Inference: compute the probability of a query/diagnosis/decision given (i.e., conditioned on) evidence/symptom/observation, i.e., $P(F | H)$

[Example from Andrew Moore]

52

Example

Statistical Inference: Compute the probability of a diagnosis, F , given symptom, H , where H = "has a headache" and F = "has flu"

That is, compute $P(F | H)$

You know that

- $P(H) = 0.1$ "one in ten people has a headache"
- $P(F) = 0.01$ "one in 100 people has flu"
- $P(H | F) = 0.9$ "90% of people who have flu have a headache"

[Example from Andrew Moore]

53

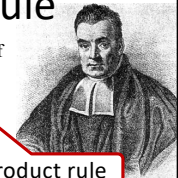
Inference with Bayes's Rule

Thomas Bayes, "Essay Towards Solving a Problem in the Doctrine of Chances," 1764

$$P(F | H) = \frac{P(F, H)}{P(H)} = \frac{P(H | F)P(F)}{P(H)}$$

Def of cond. prob.

Product rule



- $P(H) = 0.1$ "one in ten people has a headache"
- $P(F) = 0.01$ "one in 100 people has flu"
- $P(H | F) = 0.9$ "90% of people who have flu have a headache"
- $P(F | H) = (0.9 * 0.01) / 0.1 = 0.09$
- So, there's a 9% chance you have flu – much less than 90%
- But it's higher than $P(F) = 1\%$ since you have a headache

54

Bayes's Rule

- Bayes's Rule is the basis for probabilistic reasoning given a prior model of the world, $P(Q)$, and a new piece of evidence, E , Bayes's rule says how this piece of evidence decreases our ignorance about the world
- Initially, know $P(Q)$ ("prior")
- Update after knowing E ("posterior"):

$$P(Q|E) = P(Q) \frac{P(E|Q)}{P(E)}$$

55

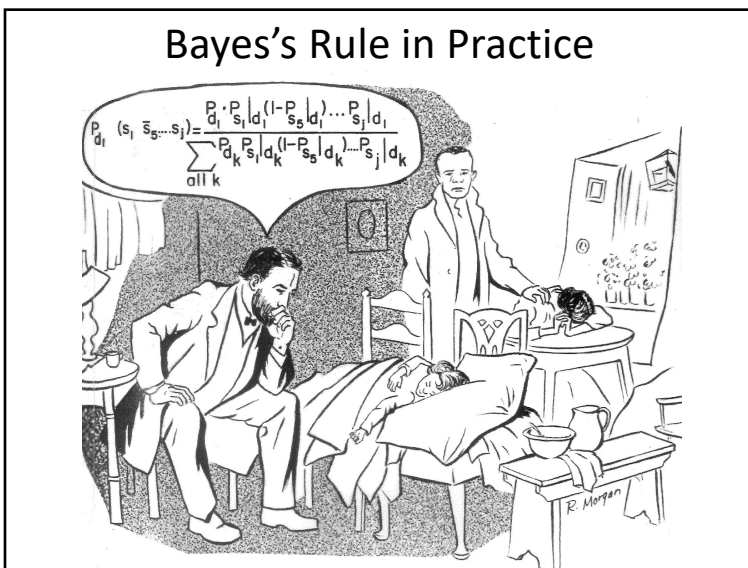
Inference with Bayes's Rule

$$P(A|B) = P(B|A)P(A) / P(B) \quad \text{Bayes's rule}$$

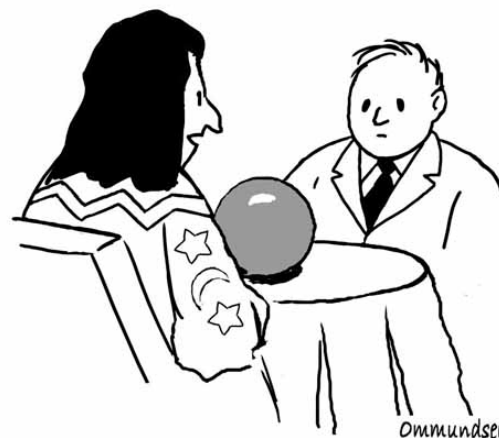
- Why do we make things this complicated?
 - Often $P(B|A)$, $P(A)$, $P(B)$ are easier to get
- Some terms:
 - **Prior:** $P(A)$: probability of A before any evidence
 - **Likelihood:** $P(B|A)$: assuming A , how likely is the evidence B
 - **Posterior:** $P(A|B)$: probability of A after knowing evidence B
 - **(Deductive) Inference:** deriving an unknown probability from known ones

56

Bayes's Rule in Practice



57



"Is this needed for a Bayesian analysis?"

58

Summary of Important Rules

- **Conditional Probability:** $P(A|B) = P(A,B)/P(B)$
- **Product rule:** $P(A,B) = P(A|B)P(B)$
- **Chain rule:** $P(A,B,C,D) = P(A|B,C,D)P(B|C,D)P(C|D)P(D)$
- **Conditionalized version of Chain rule:**

$$P(A,B|C) = P(A|B,C)P(B|C)$$
- **Bayes's rule:** $P(A|B) = P(B|A)P(A)/P(B)$
- **Conditionalized version of Bayes's rule:**

$$P(A|B,C) = P(B|A,C)P(A|C)/P(B|C)$$
- **Addition / Conditioning rule:** $P(A) = P(A,B) + P(A,\neg B)$

$$P(A) = P(A|B)P(B) + P(A|\neg B)P(\neg B)$$

59

Common Mistake

- $P(A) = 0.3$ so $P(\neg A) = 1 - P(A) = 0.7$
- $P(A|B) = 0.4$ so $P(\neg A|B) = 1 - P(A|B) = 0.6$
because $P(A|B) + P(\neg A|B) = 1$
- **but** $P(A|\neg B) \neq 0.6$ (in general)
because $P(A|B) + P(A|\neg B) \neq 1$ in general

60

Quiz

- A doctor performs a test that has 99% reliability, i.e., 99% of people who are sick test positive, and 99% of people who are healthy test negative. The doctor estimates that 1% of the population is sick.
- Question: A patient tests positive. What is the chance that the patient is sick?
- 0-25%, 25-75%, 75-95%, or 95-100%?

61

Quiz

- A doctor performs a test that has 99% reliability, i.e., 99% of people who are sick test positive, and 99% of people who are healthy test negative. The doctor estimates that 1% of the population is sick.
- Question: A patient tests positive. What is the chance that the patient is sick?
- 0-25%, 25-75%, 75-95%, or 95-100%?
- Common answer: 99%; Correct answer: 50%

62

Given:

$$P(TP | S) = 0.99$$

$$P(\neg TP | \neg S) = 0.99$$

$$P(S) = 0.01$$

TP = "tests positive"
 S = "is sick"

Query:

$$P(S | TP) = ?$$

63

$$P(TP | S) = 0.99$$

$$P(\neg TP | \neg S) = 0.99$$

$$P(S) = 0.01$$

$$P(S | TP) =$$

$$P(TP | S) P(S) / P(TP)$$

$$= (0.99)(0.01) / P(TP) = 0.0099 / P(TP)$$

$$P(\neg S | TP) = P(TP | \neg S) P(\neg S) / P(TP)$$

$$= (1 - 0.99)(1 - 0.01) / P(TP) = 0.0099 / P(TP)$$

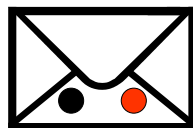
$$0.0099 / P(TP) + 0.0099 / P(TP) = 1, \text{ so } P(TP) = 0.0198$$

$$\text{So, } P(S | TP) = 0.0099 / 0.0198 = 0.5$$

64

Inference with Bayes's Rule

- In a bag there are two envelopes
 - one has a red ball (worth \$100) and a black ball
 - one has two black balls. Black balls are worth nothing



- You randomly grab an envelope, and randomly take out one ball – it's **black**
- At this point you're given the option to switch envelopes. **Should you switch or not?**

Similar to the "Monty Hall Problem"

65

Inference with Bayes's Rule

E : envelope, 1 = (R,B), 2 = (B,B)

B : the event of drawing a *black* ball

Given: $P(B | E=1) = 0.5$, $P(B | E=2) = 1$, $P(E=1) = P(E=2) = 0.5$

Query: Is $P(E=1 | B) > P(E=2 | B)$?

Use Bayes's rule: $P(E | B) = P(B | E) * P(E) / P(B)$

Conditioning rule

$$P(B) = P(B | E=1)P(E=1) + P(B | E=2)P(E=2) = (.5)(.5) + (1)(.5) = .75$$

$$P(E=1 | B) = P(B | E=1)P(E=1) / P(B) = (.5)(.5) / (.75) = 0.33$$

$$P(E=2 | B) = P(B | E=2)P(E=2) / P(B) = (1)(.5) / (.75) = 0.67$$

After seeing a black ball, the posterior probability of this envelope being #1 (thus worth \$100) is *smaller* than it being #2

Thus you should switch!

66

Another Example

- 1% of women over 40 who are tested have breast cancer. 85% of women who really *do* have breast cancer have a positive mammography test (true positive rate). 8% who do *not* have cancer will have a positive mammography (false positive rate).
- Question: A patient gets a positive mammography test. What is the chance she has breast cancer?

68

- Let Boolean random variable M mean “positive mammography test”
- Let Boolean random variable C mean “has breast cancer”
- Given:

$$P(C) = 0.01$$

$$P(M|C) = 0.85$$

$$P(M|\neg C) = 0.08$$

69

Compute the posterior probability: $P(C|M)$

70

- $P(C|M) = P(M|C)P(C)/P(M)$ by Bayes's rule
 $= (.85)(.01)/P(M)$
- $P(M) = P(M|C)P(C) + P(M|\neg C)P(\neg C)$ by the Conditioning rule
- So, $P(C|M) = .0085/[(.85)(.01) + (.08)(1-.01)]$
 $= 0.097$
- So, there is only a 9.7% chance that if you have a positive test you really have cancer!

71

Independence

Two events A, B are **independent** if the following hold:

- $P(A, B) = P(A) * P(B)$
- $P(A, \neg B) = P(A) * P(\neg B)$
- ...
- $P(A | B) = P(A)$
- $P(B | A) = P(B)$
- $P(A | \neg B) = P(A)$
- ...

74

Independence

- Independence is a kind of domain knowledge
 - Needs an understanding of **causation**
 - Very strong assumption
- Example: $P(\text{burglary}) = 0.001$ and $P(\text{earthquake}) = 0.002$
 - Let's say they are *independent*
 - The full joint probability table = ?

75

Independence

- Given: $P(B) = 0.001$, $P(E) = 0.002$, $P(B|E) = P(B)$
- The full joint probability distribution table (FJPD) is:

Burglary	Earthquake	Prob.
B	E	$= P(B)P(E)$
B	$\neg E$	
$\neg B$	E	
$\neg B$	$\neg E$	

- Need only 2 numbers to fill in entire table
- Now we can do anything, since we have the FJPD

76

Independence

- Given n independent, Boolean random variables, the FJPD has 2^n entries, but we only need n numbers (degrees of freedom) to fill in entire table
- Given n independent random variables, where each can take k values, the FJPD table has:
 - k^n entries
 - Only $n(k-1)$ numbers needed (DOFs)

77

Conditional Independence

- Random variables can be **dependent**, but **conditionally independent**
- Example: Your house has an alarm
 - Neighbor John calls when he hears the alarm
 - Neighbor Mary calls when she hears the alarm
 - Assume John and Mary don't talk to each other
- Is *JohnCall* independent of *MaryCall*?
 - **No** – If John called, it is likely the alarm went off, which increases the probability of Mary calling
 - $P(\text{MaryCall} \mid \text{JohnCall}) \neq P(\text{MaryCall})$

79

Conditional Independence

- But, if we *know* the status of the *Alarm*, *JohnCall* will **not** affect whether or not Mary calls
$$P(\text{MaryCall} \mid \text{Alarm}, \text{JohnCall}) = P(\text{MaryCall} \mid \text{Alarm})$$
- We say *JohnCall* and *MaryCall* are **conditionally independent** given *Alarm*
- In general, “**A and B are conditionally independent given C**” means:
$$P(A \mid B, C) = P(A \mid C)$$
$$P(B \mid A, C) = P(B \mid C)$$
$$P(A, B \mid C) = P(A \mid C) P(B \mid C)$$

80

Independence vs. Conditional Independence

- Say Alice and Bob each toss **separate coins**. *A* represents “Alice’s coin toss is heads” and *B* represents “Bob’s coin toss is heads”
- *A* and *B* are **independent**
- Now suppose Alice and Bob toss the **same coin**. Are *A* and *B* independent?
 - **No**. Say the coin may be biased towards heads. If *A* is heads, it will lead us to increase our belief in *B* being heads. That is, $P(B \mid A) > P(A)$

81

- Say we add a new variable, *C*: “the coin is biased towards heads”
- The values of *A* and *B* are *dependent on C*
- But if we know *for certain* the value of *C* (true or false), then any evidence about *A* cannot change our belief about *B*
- That is, $P(B \mid C) = P(B \mid A, C)$
- *A* and *B* are **conditionally independent** given *C*

82

Revisiting Earlier Example

- Let Boolean random variable M mean “positive mammography test”
- Let Boolean random variable C mean “has breast cancer”
- Given:

$$P(C) = 0.01$$

$$P(M|C) = 0.85$$

$$P(M|\neg C) = 0.08$$

85

Bayes's Rule with Multiple Evidence

- Say the same patient goes back and gets a *second* mammography and it too is positive. Now, what is the chance she has Cancer?
- Let $M1, M2$ be the 2 positive tests
- $M1$ and $M2$ are **not** independent
- Compute posterior: $P(C|M1, M2)$

86

Bayes's Rule with Multiple Evidence

- $P(C|M1, M2) = P(M1, M2|C)P(C)/P(M1, M2)$
by Bayes's rule

$$= P(M1|M2, C)P(M2|C)P(C)/P(M1, M2)$$

Conditionalized Chain rule
- $P(M1, M2) = P(M1, M2|C)P(C) +$
 $P(M1, M2|\neg C)P(\neg C)$ by Conditioning rule

$$= P(M1|M2, C)P(M2|C)P(C) +$$

$$P(M1|M2, \neg C)P(M2|\neg C)P(\neg C)$$
 by Conditionalized Chain rule

87

Cancer “causes” a positive test, so **$M1$ and $M2$ are conditionally independent given C** , so

- $P(M1|M2, C) = P(M1|C) = 0.85$
- $P(M1, M2) = P(M1|M2, C)P(M2|C)P(C) +$
 $P(M1|M2, \neg C)P(M2|\neg C)P(\neg C)$

$$= P(M1|C)P(M2|C)P(C) +$$

$$P(M1|\neg C)P(M2|\neg C)P(\neg C) \quad \text{by cond. indep.}$$

$$= (.85)(.85)(.01) + (.08)(.08)(1-.01)$$

$$= 0.01356$$

So, $P(C|M1, M2) = (.85)(.85)(.01)/.01356$

$$= 0.533 \text{ or } 53.3\%$$

88

Example

- Prior probability of having breast cancer:
 $P(C) = 0.01$
- Posterior probability of having breast cancer after 1 positive mammography:
 $P(C|M1) = 0.097$
- Posterior probability of having breast cancer after 2 positive mammographies (and cond. independence assumption):
 $P(C|M1, M2) = 0.533$

89

Bayes's Rule with Multiple Evidence

- Say the same patient goes back and gets a second mammography and it is **negative**. Now, what is the chance she has cancer?
- Let $M1$ be the positive test and $\neg M2$ be the negative test
- Compute posterior: $P(C|M1, \neg M2)$

90

Bayes's Rule with Multiple Evidence

- $P(C|M1, \neg M2) = P(M1, \neg M2|C)P(C) / P(M1, \neg M2)$
by Bayes's rule

$$= P(M1|C)P(\neg M2|C)P(C) / P(M1, \neg M2)$$

$$= (.85)(1-.85)(.01) / P(M1, \neg M2)$$
- $P(M1, \neg M2) = P(M1, \neg M2|C)P(C) + P(M1, \neg M2|\neg C)P(\neg C)$ by Conditioning rule

$$= P(M1|\neg M2, C)P(\neg M2|C)P(C) + P(M1|\neg M2, \neg C)P(\neg M2|\neg C)P(\neg C)$$

by Conditionalized Chain rule

91

Cancer "causes" a positive test, so **$M1$ and $\neg M2$ are conditionally independent given C** , so

$$\begin{aligned}
 &P(M1|\neg M2, C)P(\neg M2|C)P(C) + P(M1|\neg M2, \neg C)P(\neg M2|\neg C)P(\neg C) \\
 &= P(M1|C)P(\neg M2|C)P(C) + P(M1|\neg C)P(\neg M2|\neg C)P(\neg C) \quad \text{by cond. indep.} \\
 &= (.85)(1-.85)(.01) + (1-.08)(.08)(1-.01) \\
 &= 0.074139 \quad (= P(M1, \neg M2))
 \end{aligned}$$

$$\begin{aligned}
 \text{So, } P(C|M1, \neg M2) &= (.85)(1-.85)(.01) / .074139 \\
 &= 0.017 \text{ or } 1.7\%
 \end{aligned}$$

92

Bayes's Rule with Multiple Evidence and Conditional Independence

- Assume all evidence variables, B, C and D, are conditionally independent given the diagnosis variable, A
- $P(A|B,C,D) = P(B,C,D|A)P(A)/P(B,C,D)$
 $= \frac{P(B|A)P(C|A)P(D|A)P(A)}{P(D|B,C)P(C|B)P(B)}$

Conditionalized Chain rule + conditional independence

Chain rule

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

93

Inference Ignorance

- "Inferences about Testosterone Abuse Among Athletes," 2004
 – Mary Decker Slaney doping case
- "Justice Flunks Math," 2013
 – Amanda Knox trial in Italy

94

Naïve Bayes Classifier

- Classification problem: Find the value of class/decision/diagnosis variable Y that is most likely given evidence/measurements/attributes $X_i = v_i$
- Use Bayes's rule and conditional independence:

$$P(Y = c | X_1 = v_1, X_2 = v_2, \dots, X_n = v_n) \\ = P(Y = c)P(X_1 = v_1 | Y = c) \dots P(X_n = v_n | Y = c) / P(X_1 = v_1, \dots, X_n = v_n)$$

- Try all possible values of Y and pick the value that gives the maximum probability
- But denominator, $P(X_1 = v_1, \dots, X_n = v_n)$, is a constant for all values of Y , so it won't affect which value of Y is best

97

Naive Bayes Classifier Testing Phase

- For a given test instance defined by $X_1 = v_1, \dots, X_n = v_n$, compute

$$\operatorname{argmax}_c P(Y=c) \prod_{i=1}^n P(X_i = v_i | Y=c)$$

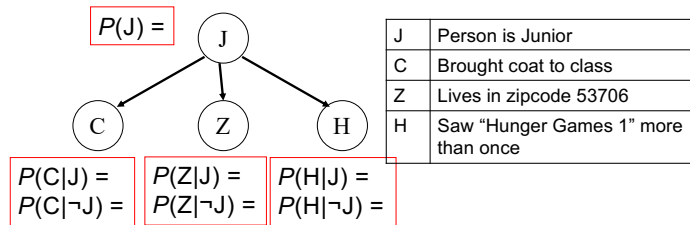
Class variable

Evidence variable

- Assumes all evidence variables are conditionally independent of each other given the class variable
- Robust because it gives the right answer as long as the correct class is more likely than all others

98

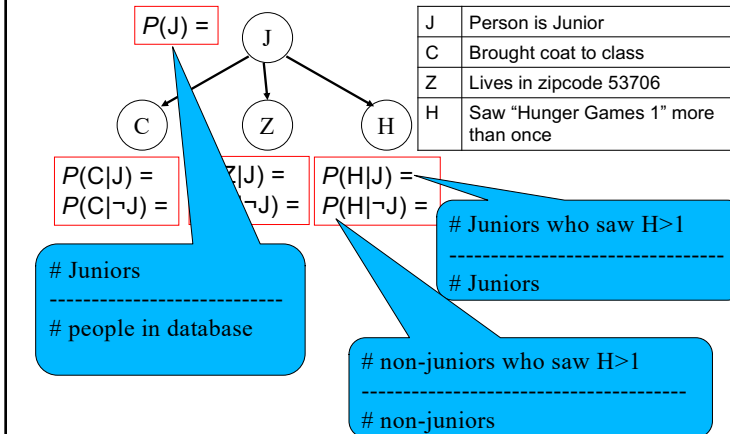
Naïve Bayes Classifier Training Phase



Compute from the Training set all the necessary Prior and Conditional probabilities

101

Naïve Bayes Classifier



102

Naïve Bayes Classifier

- Assume k classes and n evidence (i.e., attribute) variables, each with m possible values
- $k-1$ values needed for computing $P(Y=c)$
- $(m-1)k$ values needed for computing $P(X_i=v_i | Y=c)$ for each evidence variable X_i
- So, $(k-1) + n(m-1)k$ values needed instead of exponential size FJPD table

103

Naïve Bayes Classifier

- Conditional probabilities can be very, very small, so instead use logarithms to avoid underflow:

$$\arg \max_c \log P(Y = c) + \sum_{i=1}^n \log P(X_i = v_i | Y = c)$$

106

Add-1 Smoothing

- **Unseen event problem:** Training data may *not* include some cases
 - flip a coin 3 times, all heads → one-sided coin?
 - Conditional probability = 0
 - Just because a value doesn't occur in the training set doesn't mean it will never occur
- “**Add-1 Smoothing**” ensures that *every* conditional probability > 0 by pretending that you've seen each attribute's value 1 extra time

107

Add-1 Smoothing

- Compute **Conditional probabilities** as

$$P(X = v_i | Y = c) = \frac{\text{count}(X = v_i, Y = c) + 1}{\text{count}(Y = c) + m}$$

number of times
attribute X has value
 v_i in all training
instances with class c

number of training
instances with class c

- Note: $\sum_{i=1}^m P(X = v_i | Y = c) = 1$

where m = number of possible values for attribute X

108

Add-1 Smoothing

- Compute **Prior probabilities** as

$$P(Y = c) = \frac{\text{count}(Y = c) + 1}{N + k}$$

where N = size of the training set, k = number of classes

110

Laplace Smoothing

- aka **Add- δ Smoothing**
- Instead of adding 1, add δ (a positive real number)
- Compute **conditional probabilities** as

$$P(X = v_i | Y = c) = \frac{\text{count}(X = v_i, Y = c) + \delta}{\text{count}(Y = c) + \delta m}$$

where m = number of possible values for attribute X

112