



CS 5/7320

Artificial Intelligence

Uncertainty and Probabilities

AIMA Chapter 12

Slides by Michael Hahsler
based on slides by Svetlana Lazepnik
with figures from the AIMA textbook



This work is licensed under a [Creative Commons Attribution-ShareAlike 4.0 International License](https://creativecommons.org/licenses/by-sa/4.0/).

"Dice" by Steve A Johnson

Uncertainty is Bad for Agents based on Logic

Example: Catching a Flight

Let action A_t = leave for airport t minutes before flight

Question: Will A_t get me there on time?

Problems:

- Partial observability (road state, other drivers' plans, etc.)
- Noisy sensors (traffic reports)
- Uncertainty in action outcomes (flat tire, etc.)
- Complexity of modeling and predicting traffic

A purely logical approach leads to conclusions that are too weak for effective decision making:

- A_{25} will get me there on time if there is no accident on the bridge and it doesn't rain and my tires remain intact, etc., etc.
- A_{Inf} guarantees to get there in time, but who lives forever?

Making Decisions under Uncertainty

Probabilities: Suppose the agent believes the following:

$$P(A_{25} \text{ gets me there on time}) = 0.04$$

$$P(A_{90} \text{ gets me there on time}) = 0.80$$

$$P(A_{120} \text{ gets me there on time}) = 0.99$$

$$P(A_{1440} \text{ gets me there on time}) = 0.9999$$

Which action should the agent choose?

- Depends on **preferences** for missing flight vs. time spent waiting
- Encapsulated by a utility function $U(action)$

The agent should choose the action that maximizes the expected utility:

$$\operatorname{argmax}_{A_t} [P(A_t \text{ succeeds}) U(A_t \text{ succeeds}) + P(A_t \text{ fails}) U(A_t \text{ fails})]$$

- **Utility theory** is used to represent and infer preferences.
- **Decision theory** = probability theory + utility theory

Example:

Monty Hall Problem

Suppose you're on a game show, and you're given the choice of three doors: Behind one door is a car; behind the others, goats. You pick a door, say No. 3, and the host, who knows what's behind the doors, opens another door, say No. 1, which has a goat. He then says to you, "Do you want to pick door No. 2?"

Is it to your advantage to switch your choice?



Example:

Monty Hall Problem

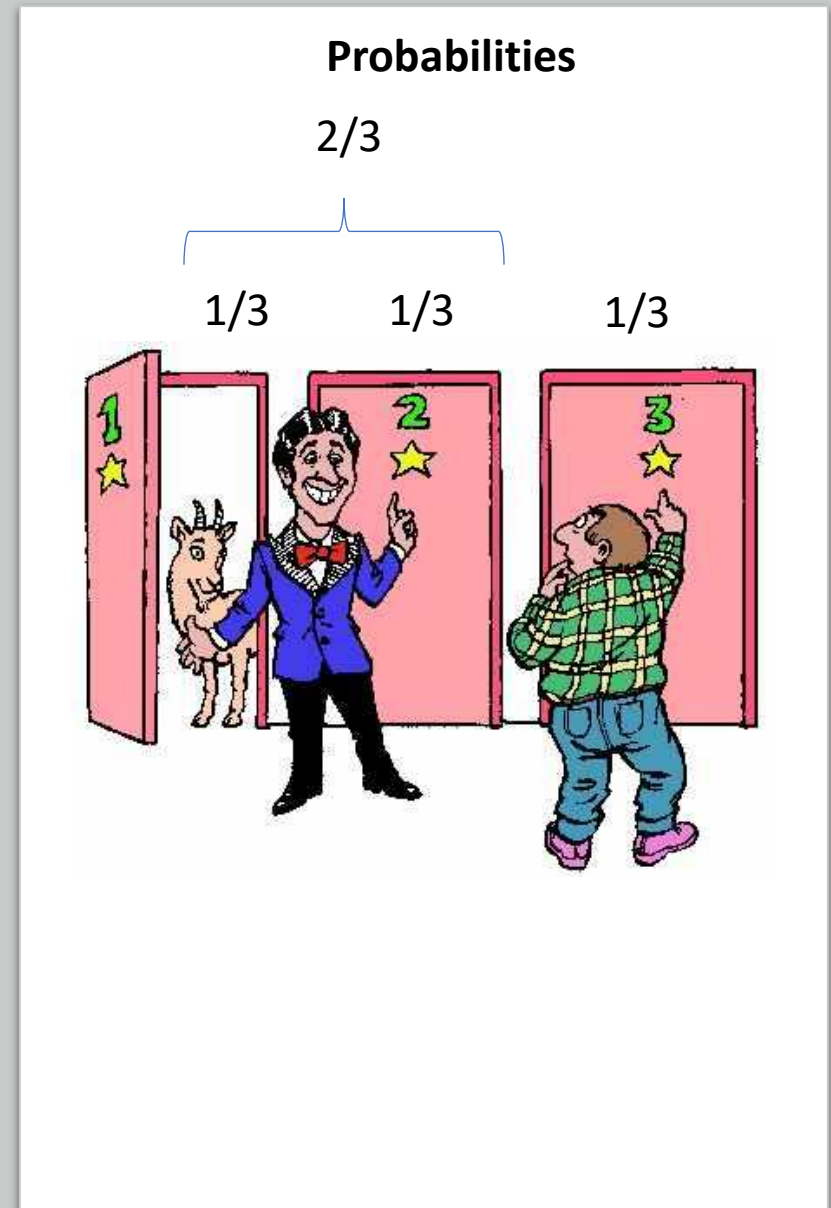
- **Assumption:** the host always opens a door with a goat!
- **Do not switch:** Initially, you have picked the correct door with probability $1/3$. If you do not switch, you have the following expected payoff:

$$(1/3) * \text{Prize} + (2/3) * 0$$

- **Switch:** Initially, there was a chance of $2/3$ that the car is behind the other two doors. The host opening a door reveals information and improves the expected payoff to:

$$(1/3) * 0 + (2/3) * \text{Prize}$$

Switch!



Sources of Uncertainty

Probabilistic assertions summarize effects of

Laziness

- failure to enumerate exceptions, qualifications, etc.

Ignorance

- lack of explicit theories, relevant facts, observability, etc.

Randomness

- Intrinsically random behavior

Example: What about a coin toss?

A Quick Review of Probability Theory

What are Probabilities?

Random variables

Events

Joint probabilities

Marginal probabilities

Conditional probabilities

Bayes' Rule

Conditional independence



What are Probabilities?

Frequentism (Empirical)

- **Probabilities are relative frequencies determined by observation.**
- For example, if we toss a coin **many times**, $P(heads)$ is the proportion of the time the coin will come up heads
- But what if we are dealing with events that only happen once? E.g., what is the probability that a Republican will win the presidency in 2024?
- **Reference class problem.** E.g., how do we define comparable elections?

Subjectivism (Bayesian Statistics)

- **Probabilities are degrees of belief updated by evidence.**
- How do we assign belief values to statements without evidence?
- How do we update our degrees of belief?
- What would make sure that agents hold consistent beliefs? E.g., The coin will land heads up and tails up at the same time.

Random variables

We describe the (uncertain) state of the world using *random variables*

- Denoted by capital letters
- **R**: *Is it raining?*
- **W**: *What's the weather?*
- **Die**: *What is the outcome of rolling two dice?*
- **V**: *What is the speed of my car (in MPH)?*

Just like variables in CSP's, random variables take on values in a *domain D*

- Domain values must be mutually exclusive and exhaustive
- **R** \in {True, False}
- **W** \in {Sunny, Cloudy, Rainy, Snow}
- **Die** \in {(1,1), (1,2), ... (6,6)}
- **V** \in [0, 200]

Events and Propositions

Probabilistic statements are defined over **events**, world states or sets of states

- *“It is raining”*
- *“The weather is either cloudy or snowy”*
- *“The sum of the two dice rolls is 11”*
- *“My car is going between 30 and 50 miles per hour”*



Events are described using

propositions:

- $R = \text{True}$
- $W = \text{“Cloudy”} \vee W = \text{“Snowy”}$
- $D \in \{(5,6), (6,5)\}$
- $30 \leq S \leq 50$

Notation:

- For random variables: $P(X = x)$, or $P(x)$ for short, is the probability of the event that random variable X has taken on the value x .
- For propositions: $P(A = \text{true})$, $P(a)$ is the probability of the set of possible worlds in which proposition A holds.

Kolmogorov's 3 Axioms of Probability

Three axioms are sufficient to define probability theory:

1. Probabilities are non-negative real numbers.
2. The probability that at least one atomic event happens is 1.
3. The probability of mutually exclusive events is additive.

This leads to important properties (A and B are sets of events):

- Numeric bound: $0 \leq P(A) \leq 1$
 - Monotonicity: if $A \subseteq B$ then $P(A) \leq P(B)$
 - Addition law: $P(A \cup B) = P(A) + P(B) - P(A \cap B)$
 - Probability of the empty set: $P(\emptyset) = 0$
 - Complement rule: $P(\neg A) = 1 - P(A)$
-
- Continuous variables need the definition of density functions.

Atomic events

- **Atomic event:** a complete specification of the state of the world, or a complete assignment of domain values to all random variables
- Atomic events are mutually exclusive and exhaustive
- E.g., if the world consists of only two Boolean variables *Cavity* and *Toothache*, then there are 4 distinct atomic events:

Cavity = false \wedge *Toothache = false*

Cavity = false \wedge *Toothache = true*

Cavity = true \wedge *Toothache = false*

Cavity = true \wedge *Toothache = true*

Joint probability distributions

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

Atomic event	P
<i>Cavity = false \wedge Toothache = false</i>	0.8
<i>Cavity = false \wedge Toothache = true</i>	0.1
<i>Cavity = true \wedge Toothache = false</i>	0.05
<i>Cavity = true \wedge Toothache = true</i>	0.05

- We use for the joint probability distribution the notation **P**(Cavity, Toothache)
- The probabilities of all possible atomic events sum to 1.

Joint probability distributions

- Suppose we have a joint *distribution* $P(X_1, X_2, \dots, X_n)$ of n random variables with domain sizes d
 - **The size of the probability table is $d \times d \times \dots \times d = d^n$**
 - Impossible to write out completely for all distributions with many variables! We will come back to this problem when we talk about independence.
- Notation:
 - $P(x)$, $P(X=x)$ is the **probability** that random variable X takes on value x
 - $P(X)$ is the **distribution of probabilities** for all possible values of X . Often we are lazy or forget to make P bold.

Marginal probability distributions

- Sometimes we are only interested in one variable. This is called the *marginal distribution* $P(Y)$

P(Cavity, Toothache)	
$Cavity = false \wedge Toothache = false$	0.8
$Cavity = false \wedge Toothache = true$	0.1
$Cavity = true \wedge Toothache = false$	0.05
$Cavity = true \wedge Toothache = true$	0.05

Marginal
Prob. Distr.

P(Cavity)	
$Cavity = false$?
$Cavity = true$?

P(Toothache)	
$Toothache = false$?
$Toothache = true$?

Marginal probability distributions

- Suppose we have the joint distribution $P(X,Y)$ and we want to find the *marginal distribution* $P(Y)$

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

- **General rule:** to find $P(X = x)$, sum the probabilities of all atomic events where $X = x$. This is called “summing out” or marginalization.

Marginal probability distributions

- Suppose we have the joint distribution $P(X,Y)$ and we want to find the *marginal distribution* $P(Y)$

P(Cavity, Toothache)	
$Cavity = false \wedge Toothache = false$	0.8
$Cavity = false \wedge Toothache = true$	0.1
$Cavity = true \wedge Toothache = false$	0.05
$Cavity = true \wedge Toothache = true$	0.05

Marginal Prob. Distr.	P(Cavity)	
	$Cavity = false$	$0.8+0.1 = 0.9$
	$Cavity = true$	$0.05+0.05=0.1$

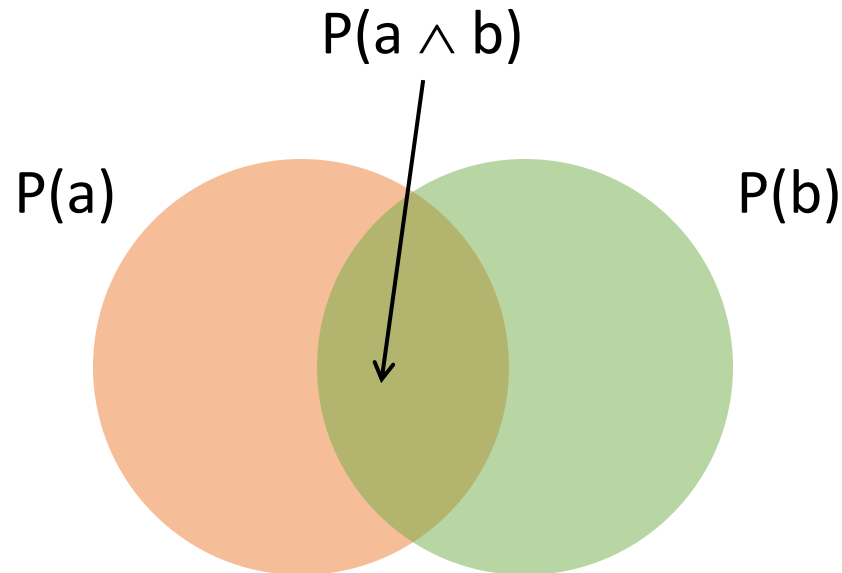
P(Toothache)	
$Toothache = false$	$0.8+0.05= 0.85$
$Toothache = true$	$0.1+0.05= 0.15$

Conditional probability

- Probability of cavity given toothache:

$$P(\text{Cavity} = \text{true} \mid \text{Toothache} = \text{true})$$

- For any two events a and b , $P(a|b) = \frac{P(a \wedge b)}{P(b)} = \frac{P(a,b)}{P(b)}$



Conditional probability

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

Joint Prob. Distr.	P(Cavity, Toothache)	
	<i>Cavity = false</i> \wedge <i>Toothache = false</i>	0.8
	<i>Cavity = false</i> \wedge <i>Toothache = true</i>	0.1
	<i>Cavity = true</i> \wedge <i>Toothache = false</i>	0.05
	<i>Cavity = true</i> \wedge <i>Toothache = true</i>	0.05

Marginal Prob. Distr.	P(Cavity)	
	<i>Cavity = false</i>	0.9
	<i>Cavity = true</i>	0.1

P(Toothache)	
<i>Toothache = false</i>	0.85
<i>Toothache = true</i>	0.15

- What is $P(\text{Cavity} = \text{true} \mid \text{Toothache} = \text{false})$?
 $0.05 / 0.85 = 0.059$
- What is $P(\text{Cavity} = \text{false} \mid \text{Toothache} = \text{true})$?
 $0.1 / 0.15 = 0.667$

Conditional distributions

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

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

A conditional distribution is a distribution over the values of one variable given fixed values of other variables

P(Cavity Toothache = true)	
<i>Cavity = false</i>	0.667
<i>Cavity = true</i>	0.333

P(Cavity Toothache = false)	
<i>Cavity = false</i>	0.941
<i>Cavity = true</i>	0.059

P(Toothache Cavity = true)	
<i>Toothache = false</i>	0.5
<i>Toothache = true</i>	0.5

P(Toothache Cavity = false)	
<i>Toothache = false</i>	0.889
<i>Toothache = true</i>	0.111

Normalization trick

- To get the whole conditional distribution $P(X | y)$ at once, select all entries in the joint distribution matching $Y = y$ and renormalize them to sum to one

$P(\text{Cavity}, \text{Toothache})$	
$\text{Cavity} = \text{false} \wedge \text{Toothache} = \text{false}$	0.8
$\text{Cavity} = \text{false} \wedge \text{Toothache} = \text{true}$	0.1
$\text{Cavity} = \text{true} \wedge \text{Toothache} = \text{false}$	0.05
$\text{Cavity} = \text{true} \wedge \text{Toothache} = \text{true}$	0.05



Select $P(X, y)$

$\text{Toothache}, \text{Cavity} = \text{false}$	
$\text{Toothache} = \text{false}$	0.8
$\text{Toothache} = \text{true}$	0.1



Sum is $P(y) = 0.9$



Renormalize sum to 1 (= divide by $P(y)$)

$P(\text{Toothache} \text{Cavity} = \text{false})$	
$\text{Toothache} = \text{false}$	0.889
$\text{Toothache} = \text{true}$	0.111

Equivalent to

$$P(X | y) = \alpha P(X, y)$$

with $\alpha = 1/P(y)$

Bayes' Rule

- The **product rule** gives us two ways to factor a joint distribution:

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

Posterior Prob.

Prior Prob.

- Therefore,
$$P(A|B) = \frac{P(B|A) P(A)}{P(B)}$$

- Why is this useful?

- Can get *diagnostic probability* $P(\text{cavity} \mid \text{toothache})$ from *causal probability* $P(\text{toothache} \mid \text{cavity})$
- We can update our beliefs based on evidence.
- Important tool for probabilistic inference .

Rev. Thomas Bayes
(1702-1761)

Example: Getting Married in the Desert

Marie is getting married tomorrow, at an outdoor ceremony in the desert. In recent years, it has rained only 5 days each year ($5/365 = 0.014$). Unfortunately, the weatherman has predicted rain for tomorrow. When it actually rains, the weatherman correctly forecasts rain 90% of the time. When it doesn't rain, he incorrectly forecasts rain 10% of the time. What is the probability that it will rain on Marie's wedding?

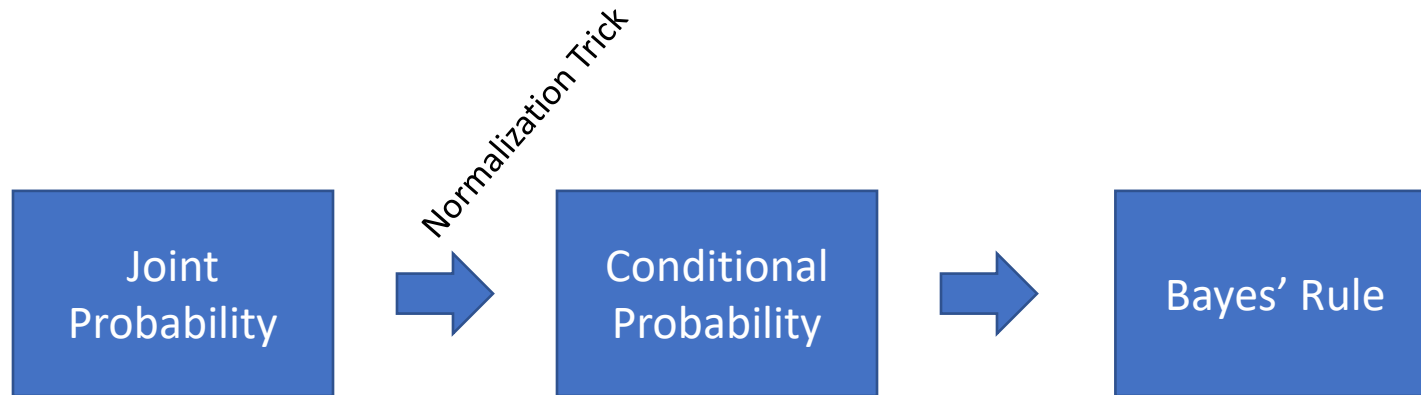
Example: Getting Married in the Desert

Marie is getting married tomorrow, at an outdoor ceremony in the desert. In recent years, it has rained only 5 days each year ($5/365 = 0.014$). Unfortunately, the **weatherman has predicted rain** for tomorrow. When it actually rains, the weatherman correctly forecasts rain 90% of the time. When it doesn't rain, he incorrectly forecasts rain 10% of the time. What is the **probability that it will rain** on Marie's wedding?

$$\begin{aligned} P(\text{Rain}|\text{Predict}) &= \frac{P(\text{Predict}|\text{Rain})P(\text{Rain})}{P(\text{Predict})} \\ &= \frac{P(\text{Predict}|\text{Rain})P(\text{Rain})}{P(\text{Predict}|\text{Rain})P(\text{Rain}) + P(\text{Predict}|\neg\text{Rain})P(\neg\text{Rain})} \\ &= \frac{0.9 * 0.014}{0.9 * 0.014 + 0.1 * 0.986} = 0.111 \end{aligned}$$

The weather forecast updates our belief from 0.014 to 0.111

Approach



- **Problem:** Joint probability table is typically too large! For n random variables with a domain size of d we have a table of size $O(d^n)$. This is a problem for
 - storing the table and
 - Estimating the probabilities.
- **Solution:** Decomposition of joint probability distributions using **independence** and conditional independence between events. A large table can be broken into several much smaller tables.

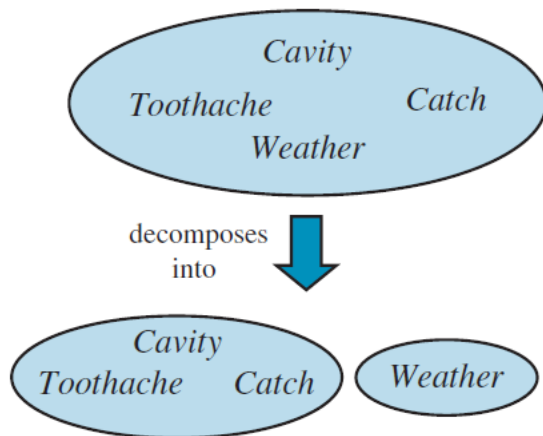
Independence between Events

- Two events a and b are **independent** if and only if

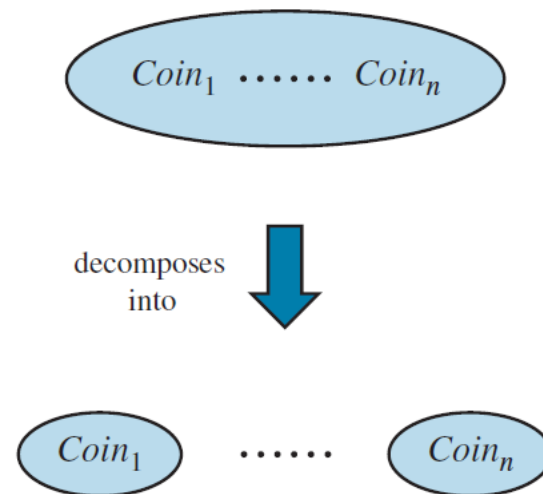
$$P(a \wedge b) = P(a) P(b)$$

- This is equivalent to $P(a \mid b) = P(a)$ and $P(b \mid a) = P(b)$
- Independence is an important simplifying assumption for modeling, e.g., *Cavity* and *Weather* can be assumed to be independent

$$P(\text{Cavity} \mid \text{Weather}) = P(\text{Cavity})$$



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



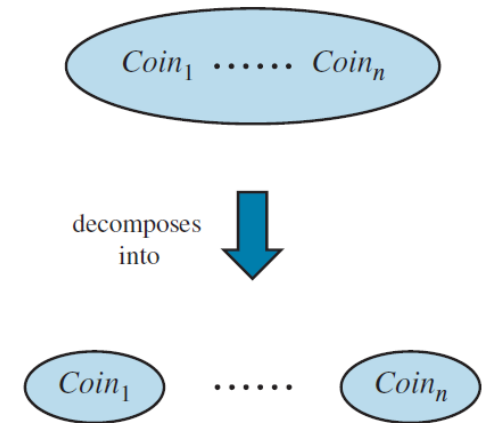
$$P(H, H, T, \dots, T) = P(H)P(H)P(T)\dots P(T)$$

Decomposition of the Joint Probability Distribution

- **Independence:** The joint probability can be decomposed into

$$\begin{aligned} &P(Coin_1, \dots, Coin_n) \\ &= P(Coin_1) \times \dots \times P(Coin_n) = \prod_i P(Coin_i) \end{aligned}$$

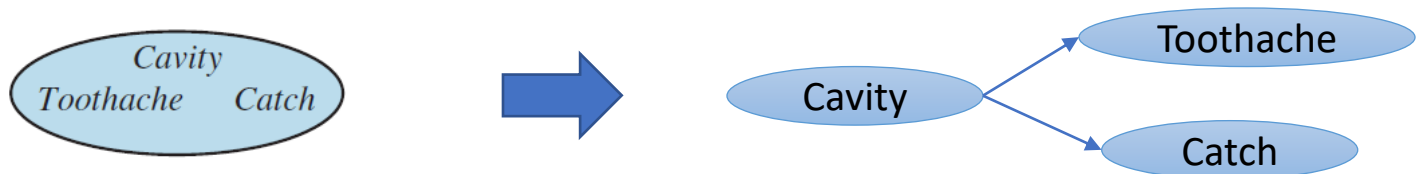
- We need for each coin one parameter (chance of getting H).
- Independence reduces the numbers needed to specify the joint distribution from $2^n - 1$ to n .
- Note: If we have identical (iid) coins, then we even only need 2 numbers, probability of H and number of coins.



Conditional Independence

- **Conditional independence:** a and b are *conditionally independent* given c (i.e., if we know c) iff

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



- If the patient has a cavity, the probability that the probe catches in it does not depend on whether he/she has a toothache

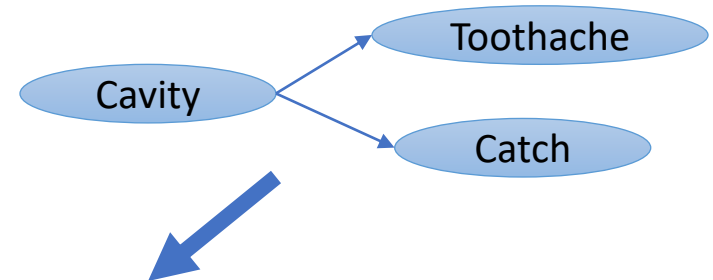
$$P(\text{Catch} \mid \text{Toothache}, \text{Cavity}) = P(\text{Catch} \mid \text{Cavity})$$

- Therefore, *Catch* is **conditionally independent** of *Toothache* given *Cavity*
- Likewise, *Toothache* is conditionally independent of *Catch* given *Cavity*

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

Decomposition of the Joint Probability Distribution

- **Conditional independence**
using the chain rule:



$$\begin{aligned} P(\text{Toothache}, \text{Catch}, \text{Cavity}) &= \\ P(\text{Cavity}) P(\text{Catch} \mid \text{Cavity}) P(\text{Toothache} \mid \cancel{\text{Catch}}, \text{Cavity}) &= \\ P(\text{Cavity}) P(\text{Catch} \mid \text{Cavity}) P(\text{Toothache} \mid \text{Cavity}) \end{aligned}$$

- The full joint probability distribution needs $2^3 - 1 = 7$ independent numbers (-1 because the 2^3 numbers have to sum up to 1).
- Conditional independence reduces this to: $1+2+2=5$
- In many practical applications, conditional independence reduces the space requirements from $O(2^n)$ to $O(n)$.

Chain rule: Example for 4 variables.

$$\begin{aligned} P(X_4, X_3, X_2, X_1) &= P(X_4 \mid X_3, X_2, X_1) \cdot P(X_3, X_2, X_1) \\ &= P(X_4 \mid X_3, X_2, X_1) \cdot P(X_3 \mid X_2, X_1) \cdot P(X_2, X_1) \\ &= P(X_4 \mid X_3, X_2, X_1) \cdot P(X_3 \mid X_2, X_1) \cdot P(X_2 \mid X_1) \cdot P(X_1) \end{aligned}$$



Bayesian Decision Making

Making Decisions under Uncertainty based on Evidence

Probabilities and Rationality

Why should a rational agent hold beliefs that are consistent with axioms of probability?

- **De Finetti (1931):** If an agent has some degree of belief in proposition A, he/she should be able to decide whether to accept a bet for/against A.
- **Example:** If an agent believes that $P(A) = 0.4$, should he/she agree to bet \$4 that A will occur against \$6 that A will not occur?

The utility would be $(0.4 \times 10 + 0.6 \times 0) - 4 = 0$. Not having $P(A) = 1 - P(\neg A)$ would be a problem!

- **Theorem:** An agent who holds beliefs inconsistent with axioms of probability can be tricked into accepting a combination of bets that guarantees that the agent loses money.

Probabilistic inference

Suppose the agent has to decide about the value of an unobserved *query variable* X given some observed *evidence* $E = e$ and we assume X causes E .

- Examples:

- $x \in \{\text{spam}, \text{not spam}\}$, e = email message

- $x \in \{\text{zebra}, \text{giraffe}, \text{hippo}\}$, e = image features

Bayes' decision theory

- The agent has a **loss function**, which is 0 if the value of X is guessed correctly and 1 otherwise.
- The estimate of X that minimizes the *expected loss* is the one that has the greatest posterior probability $P(X = x \mid E = e) = P(x|e)$.
- This is called the **MAP** (maximum a posteriori) decision.

MAP: Maximum A Posteriori Decision

Value of x that has the highest (maximum) posterior probability given the evidence e

Posterior Prob.

$$\begin{aligned} x^* &= \operatorname{argmax}_x \overbrace{P(x|e)}^{\text{Posterior Prob.}} = \operatorname{argmax}_x \frac{P(e|x)P(x)}{P(e)} \\ &\propto \operatorname{argmax}_x P(e|x) \underbrace{P(x)}_{\text{Prior Prob.}} \end{aligned}$$

For comparison: the maximum likelihood decision ignores $P(X)$

$$x^* = \operatorname{argmax}_x \underbrace{P(e|x)}_{\text{likelihood}}$$

MAP: Example

Value of x that has the highest (maximum) posterior probability given the evidence e .

$x \in \{zebra, dog, cow\}, e = \text{stripes}$

Posterior Prob.

$$\begin{aligned} x^* = \operatorname{argmax}_x \overbrace{P(x|e)} &= \operatorname{argmax}_x \frac{P(\text{stripes}|x)P(x)}{P(\text{stripes})} \\ &\propto \operatorname{argmax}_x \underbrace{P(\text{stripes}|x)}_{\text{likelihoods}} \underbrace{P(x)}_{\text{Prior Prob.}} \end{aligned}$$

$P(\text{stripes}|\text{zebra})$ is obviously the highest, but it also depends on the prior $P(\text{zebra})$, the chance that we see a zebra (which would be higher in a zoo than in downtown Dallas).

Bayes Classifier

- Suppose we have many different types of observations (evidence, symptoms, features) F_1, \dots, F_n that we want to use to obtain evidence about an underlying hypothesis H
- MAP decision involves estimating

$$P(H|F_1, \dots, F_n) \propto P(F_1, \dots, F_n|H)P(H)$$

- If each feature can take on k values, how many entries are in the joint probability table $P(F_1, \dots, F_n|H)$?

Naïve Bayes model

- Suppose we have many different types of observations (evidence, symptoms, features) F_1, \dots, F_n that we want to use to obtain evidence about an underlying hypothesis H
- MAP decision involves estimating

$$P(H|F_1, \dots, F_n) \propto P(F_1, \dots, F_n|H)P(H)$$

- We can make the simplifying assumption that the different **features are conditionally independent given the hypothesis** reduces the joint probability distribution table to size $k \times n$:

$$P(F_1, \dots, F_n|H) = \prod_{i=1}^n P(F_i|H)$$

Example: Naïve Bayes Spam Filter

- The hypothesis can be *spam* or $\neg\text{spam}$ and the evidence is the message.
- **MAP decision:** to minimize the probability of error, we should classify a message as spam if

$$P(\text{spam} \mid \text{message}) > P(\neg\text{spam} \mid \text{message})$$



Dear Sir.

First, I must solicit your confidence in this transaction, this is by virtue of its nature as being utterly confidential and top secret. ...



TO BE REMOVED FROM FUTURE MAILINGS, SIMPLY REPLY TO THIS MESSAGE AND PUT "REMOVE" IN THE SUBJECT.

99 MILLION EMAIL ADDRESSES
FOR ONLY \$99



Ok, I know this is blatantly OT but I'm beginning to go insane. Had an old Dell Dimension XPS sitting in the corner and decided to put it to use, I know it was working pre being stuck in the corner, but when I plugged it in, hit the power nothing happened.

**How do we represent
the messages?**

Natural Language Processing: Bag of Words

Represent a document as binary vector (w_1, \dots, w_n) . Each element represents the event that word w_i is present ($w_i = 1$) or not ($w_i = 0$).

Simplifications:

- The order of the words in the message is ignored.
- How often a word is repeated is ignored.



Dear Sir.

First, I must solicit your confidence in this transaction, this is by virtue of its nature as being utterly confidential and top secret. ...



TO BE REMOVED FROM FUTURE MAILINGS, SIMPLY REPLY TO THIS MESSAGE AND PUT "REMOVE" IN THE SUBJECT.

99 MILLION EMAIL ADDRESSES
FOR ONLY \$99



Ok, I know this is blatantly OT but I'm beginning to go insane. Had an old Dell Dimension XPS sitting in the corner and decided to put it to use, I know it was working pre being stuck in the corner, but when I plugged it in, hit the power nothing happened.

Naïve Bayes Spam Filter

- If we assume that each word is conditionally independent of the others given message class (spam or not spam), then we can use a naïve Bayes classifier.

$$P(message|spam) = P(w_1, \dots, w_n|spam) = \prod_{i=1}^n P(w_i|spam)$$

$$\underbrace{P(spam|w_1, \dots, w_n)}_{\text{posterior}} \propto \underbrace{P(spam)}_{\text{prior}} \underbrace{\prod_{i=1}^n P(w_i|spam)}_{\substack{\text{Evidence} \\ \text{(presents and absence of words)}}}$$

Parameter estimation

In order to classify a message, we need to know

1. the prior $P(\text{spam})$ and
2. the likelihoods $P(\text{word} = 1 \mid \text{spam})$, $P(\text{word} = 0 \mid \text{spam})$, $P(\text{word} = 1 \mid \neg\text{spam})$ and $P(\text{word} = 0 \mid \neg\text{spam})$

These are the *parameters* of the probabilistic model:

prior

spam:	0.33
\neg spam:	0.67

$P(\text{word} = 1 \mid \text{spam})$

the :	0.0156
to :	0.0153
and :	0.0115
of :	0.0095
you :	0.0093
a :	0.0086
with:	0.0080
from:	0.0075
...	

$P(\text{word} = 1 \mid \neg\text{spam})$

the :	0.0210
to :	0.0133
of :	0.0119
2002:	0.0110
with:	0.0108
from:	0.0107
and :	0.0105
a :	0.0100
...	

$$P(\text{word} = 0 \mid \text{spam}) = 1 - P(\text{word} = 1 \mid \text{spam})$$

$$P(\text{word} = 0 \mid \neg\text{spam}) = 1 - P(\text{word} = 1 \mid \neg\text{spam})$$

Parameter estimation: Prior

How do we obtain the prior $P(\text{spam})$?

Empirically: use training data

$$P(\text{spam}) = \frac{\text{\textit{\# of spam messages}}}{\text{\textit{total \# of messages}}}$$

$$P(\neg\text{spam}) = 1 - P(\text{spam})$$

Parameter estimation: Likelihoods

How do we obtain the likelihoods $P(\text{word} = 1 \mid \text{spam})$ and $P(\text{word} = 1 \mid \neg\text{spam})$?

Empirically: use training data

$$P(\text{word} = 1 \mid \text{spam}) = \frac{\text{\textit{\# of spam messages that contain the word}}}{\text{\textit{total \# of spam messages}}}$$

Note: $P(\text{word} \mid \text{spam})$ is the likelihood and the equation above is the *maximum likelihood* (ML) estimate. The estimate that maximizes probability of the data (words) given the parameter (class):

$$L = \prod_{d=1}^D \prod_{i=1}^{n_d} P(w_{d,i} \mid \text{class}_{d,i})$$

d : index of training document

i : index of a word

Parameter estimation: Smoothing

- **Problem:** What happens with words that we have never seen or seen only a few times?
- **Laplacian smoothing:** add one to each count

$$P(\text{word} = 1 \mid \text{spam}) = \frac{\text{\textit{\# of spam messages that contain the word}} + 1}{\text{\textit{total \# of spam messages}} + \text{\textit{\# of classes}}}$$

Note: This is actually a Bayesian estimate with +1 and # of classes (2 for spam/not spam) representing the prior probability.

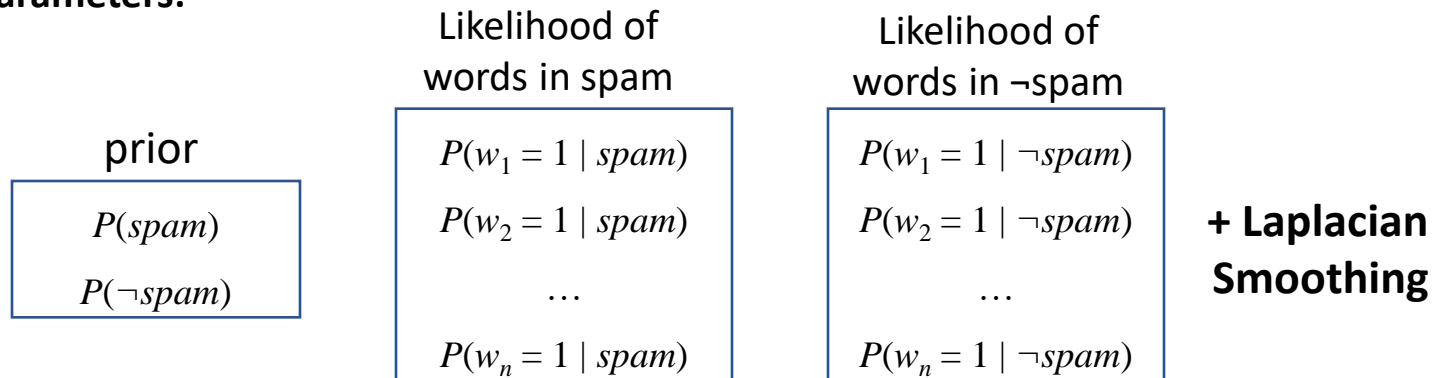
Summary of model and parameters

- **Naïve Bayes model:**

$$P(spam|message) \propto P(spam) \prod_{i=1}^n P(w_i|spam) = score_{spam}(message)$$

$$P(\neg spam|message) \propto P(\neg spam) \prod_{i=1}^n P(w_i|\neg spam) = score_{\neg spam}(message)$$

- **Model parameters:**



Likelihood of words not in spam (or $\neg spam$) can be calculated as
 $P(w_i = 0 | spam) = 1 - P(w_i = 1 | spam)$

- **Decision:** Spam if

$$P(spam | message) > P(\neg spam | message)$$

equivalent to

$$score_{spam}(message) > score_{\neg spam}(message)$$

Bayesian decision making: Summary

- Suppose the agent has to decide about the value of an unobserved *query variable* X based on the values of an observed *evidence variable* E
- **Inference problem:** given some evidence $E = e$, what is the posterior probability $P(X | e)$?
- **Learning problem:** estimate the parameters of the probabilistic model $P(X | E)$ and $P(X)$ given a *training sample* $\{(x_1, e_1), \dots, (x_n, e_n)\}$
- Learning from data is the goal of **Machine Learning**.