

# Mathematics Review

CS-585

**Natural Language Processing**

Sonjia Waxmonsky

Slides based in part on material from:

- *Artificial Intelligence: A Modern Approach, 2nd Edition*  
Russell & Norvig (Prentice-Hall: 2003)
- Slides by Patrick Nichols (MIT), Derrick Higgins (IIT)

---

# PROBABILITY THEORY REVIEW

# Probability: Vocabulary

---

*Some concepts we will cover today:*

- *Complement*
- *Conditional probability*
- *Prior probability*
- *Posterior probability*
- *Random variable (binary vs multi-valued)*
- *Independent variables*
- *Conditional independence*
- *Chain rule*
- *Atomic event*
- *Bayes Rule*

# Probability and NLP

---

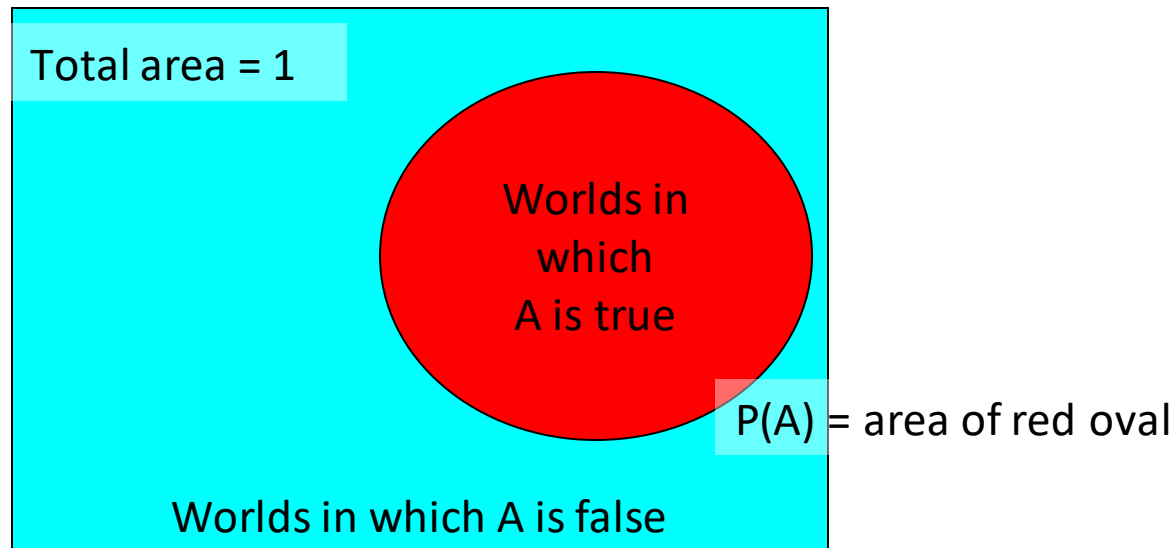
Why do we need probability for modern NLP?

- Based on what we can observe, what is the *most likely* label or outcome of multiple possibilities?
  - *What is the best interpretation of "Time flies like an arrow"?* [Parsing, part-of-speech tagging]
  - *How many stars would the writer of this review give to this movie?* [Sentiment analysis]
  - *What is the best response to "Who is the vice president"?* [Question answering]

# Probability: Intuitive

- $P(A)$  denotes “fraction of possible worlds (given what I know) in which  $A$  is true”

Event Space of all possible worlds



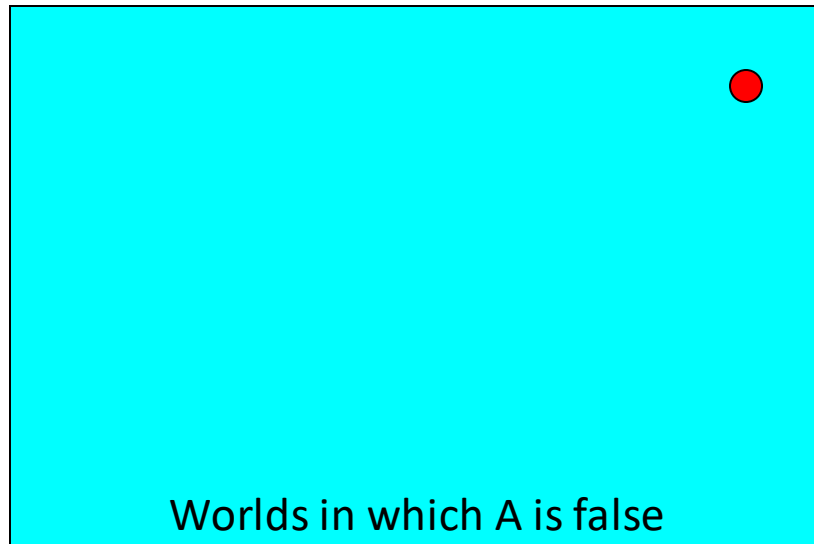
# Probability: Axioms

---

- $0 \leq P(A) \leq 1$
- $P(\text{true}) = 1$
- $P(\text{false}) = 0$
- $P(A \text{ or } B) = P(A) + P(B) - P(A \ \& \ B)$

# Probability: Axioms

- $0 \leq P(A) \leq 1$
- $P(\text{true}) = 1$
- $P(\text{false}) = 0$
- $P(A \text{ or } B) = P(A) + P(B) - P(A \ \& \ B)$

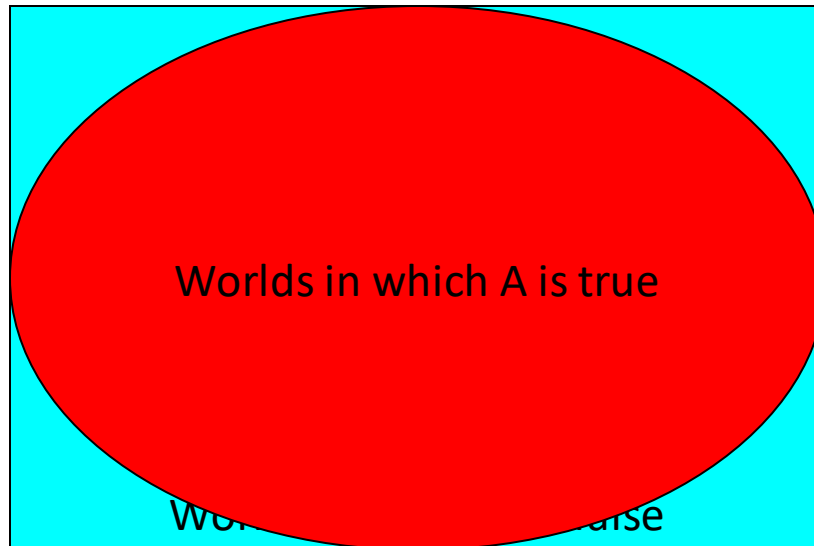


Red oval can't get smaller than 0

Area of 0 means that A is true in  
**no** possible worlds...

# Probability: Axioms

- $0 \leq P(A) \leq 1$
- $P(\text{true}) = 1$
- $P(\text{false}) = 0$
- $P(A \text{ or } B) = P(A) + P(B) - P(A \ \& \ B)$



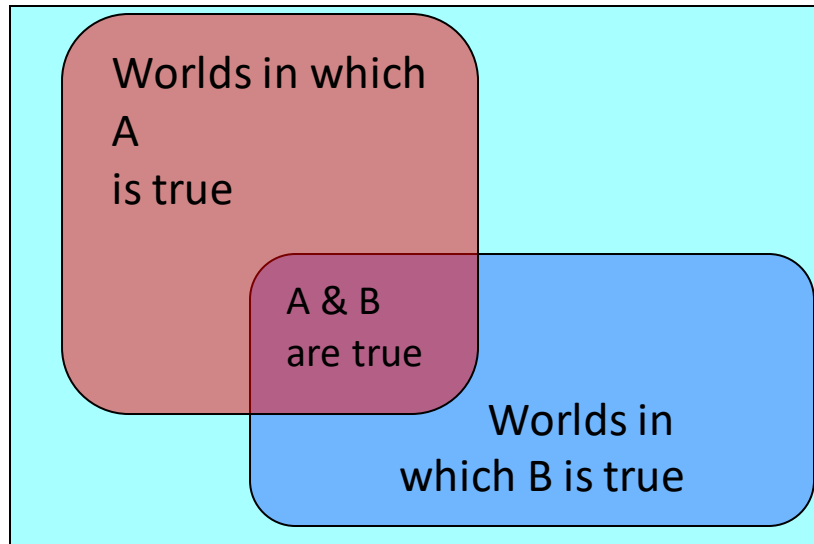
Red oval can't get larger than 1

Area of 1 means that A is true in **all** possible worlds...



# Probability: Axioms

- $0 \leq P(A) \leq 1$
- $P(\text{true}) = 1$
- $P(\text{false}) = 0$
- $P(A \text{ or } B) = P(A) + P(B) - P(A \ \& \ B)$



Size of union is sum of sizes  
minus size of intersection

# Some Provable Facts

## Axioms:

- $0 \leq P(A) \leq 1$
- $P(\text{true}) = 1$
- $P(\text{false}) = 0$
- $P(A \text{ or } B) = P(A) + P(B) - P(A \ \& \ B)$

We can show that:

- $P(\sim A) = P(\text{not } A) = 1 - P(A)$

And furthermore:

- $P(A) = P(A \ \& \ B) + P(A \ \& \ \sim B)$

Here  $P(\sim A)$  is the *complement* of A

# Multivalued Random Variables

---

Suppose  $A$  can take on more than 2 values

**Example:**

Part of Speech ( $POS$ ): {noun, verb, adjective, adverb}

Call  $A$  a **random variable with arity  $k$**  if  $A$  can take on one of  $k$  different values in some set  $\{v_1, v_2, \dots, v_k\}$

Thus:

- $P(A=v_i \ \& \ A=v_j) = 0 \quad \text{if } i \neq j$
- $P(A=v_1 \text{ or } A=v_2 \text{ or } \dots \text{ or } A=v_k) = 1$

# Easy Facts About Multivalued RVs

## Axioms:

- $0 \leq P(A) \leq 1$ ;  $P(\text{true}) = 1$ ;  $P(\text{false}) = 0$
- $P(A \text{ or } B) = P(A) + P(B) - P(A \& B)$

## Recall:

- $P(A=v_i \& A=v_j) = 0$  if  $i \neq j$ ;  $P(A=v_1 \text{ or } A=v_2 \text{ or } \dots \text{ or } A=v_k) = 1$

- We can show that:

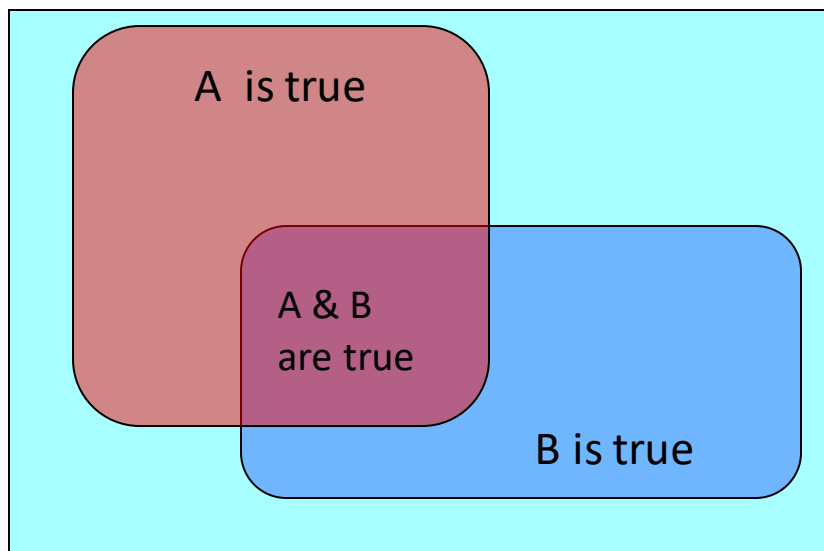
$$P(A = v_1 \vee A = v_2 \vee \dots \vee A = v_i) = \sum_{j=1}^i P(A = v_j)$$

- And therefore:

$$P(A = v_1 \vee \dots \vee A = v_k) = \sum_{j=1}^k P(A = v_j)$$

# Conditional Probability

- $P(A | B)$  = “probability of  $A$  **given**  $B$ ” = fraction of possible worlds with  $B$  true that also have  $A$  true



$$P(\text{Headache}) = 0.1$$

$$P(\text{Flu}) = 0.02$$

$$P(\text{Headache} | \text{Flu}) = 0.5$$

“Headaches are rare, Flu is much rarer, but if you have the Flu, you have a 50-50 chance of having a headache.”

# Conditional Probability

---

- *Formal definition:*

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

- *This gives us:*

$$P(A \wedge B) = P(A | B)P(B)$$

# Chain Rule

---

From Conditional Probability:

$$P(A \wedge B) = P(A/B) P(B)$$

Thus, we have:

$$\begin{aligned} P(A \wedge B \wedge C) &= P(A/B \wedge C) P(B \wedge C) \\ &= P(A/B \wedge C) P(B/C) P(C) \end{aligned}$$

Generalizing:

$$P(A1 \wedge A2 \wedge \dots \wedge An) = P(A1/A2 \wedge \dots \wedge An) P(A2/A3 \wedge \dots \wedge An) P(An)$$

# Atomic Events

---

- **Atomic event**: A **complete** specification of the state of the world about which the agent is uncertain

E.g., if the world consists of only two Boolean variables *Cavity* and *Toothache*, then there are 4 distinct atomic events:

*Cavity = false & Toothache = false*

*Cavity = false & Toothache = true*

*Cavity = true & Toothache = false*

*Cavity = true & Toothache = true*

- Atomic events are mutually exclusive and exhaustive



# Prior probability

- **Prior** or **unconditional probabilities** of propositions

e.g.,  $P(\text{Cavity} = \text{true}) = 0.1$  and  $P(\text{Weather} = \text{sunny}) = 0.72$  correspond to belief prior to arrival of any (new) evidence

- **Probability distribution** gives values for all possible assignments:

$P(\text{Weather}) = \langle 0.72, 0.1, 0.08, 0.1 \rangle$  (**normalized**, i.e., sums to 1)

- **Joint probability distribution** for a set of random variables gives the probability of every atomic event on those random variables

$P(\text{Weather}, \text{Cavity})$  = a  $4 \times 2$  matrix of values:

| <i>Weather</i> =      | sunny | rainy | cloudy | snow |
|-----------------------|-------|-------|--------|------|
| <i>Cavity</i> = true  | 0.144 | 0.02  | 0.016  | 0.02 |
| <i>Cavity</i> = false | 0.576 | 0.08  | 0.064  | 0.08 |

- Every question about a domain can be answered by the joint distribution

# Inference

---

- Generally: Given some information about the probability distribution, determine the probability of some proposition  $\phi$
- $\phi = \textit{Cavity}$
- $\phi = \textit{Cavity} \ \& \ \textit{Toothache}$
- $\phi = \sim \textit{Study} \ \& \ (\textit{GoodGrade} \ \mathbf{or} \ \textit{GoodJob})$

# Inference by enumeration

- Start with the joint probability distribution:

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

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

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

# Inference by enumeration

- Start with the joint probability distribution:

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

- For any proposition  $\phi$ , sum the atomic events where it is true:  $P(\phi) = \sum_{\omega: \omega \models \phi} P(\omega)$
- $P(\text{toothache}) = 0.108 + 0.012 + 0.016 + 0.064 = 0.2$

# Inference by enumeration

- Start with the joint probability distribution:

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

- For any proposition  $\phi$ , sum the atomic events where it is true:  $P(\phi) = \sum_{\omega: \omega \models \phi} P(\omega)$
- $P(\text{toothache or cavity}) =$   
 $0.108 + 0.012 + 0.016 + 0.064 + 0.072 + 0.008 = 0.28$

# Inference by enumeration

- Start with the joint probability distribution:

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

- Can also compute conditional probabilities:

$$P(\sim\text{cavity} \mid \text{toothache}) \quad ?$$

# Inference by enumeration

- Start with the joint probability distribution:

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

- Can also compute conditional probabilities:

$$\begin{aligned} P(\sim\text{cavity} \mid \text{toothache}) &= \frac{P(\sim\text{cavity} \ \& \ \text{toothache})}{P(\text{toothache})} \\ &= \frac{0.016 + 0.064}{0.108 + 0.012 + 0.016 + 0.064} \\ &= 0.4 \end{aligned}$$

# Independence

---

Two boolean random variables A and B are said to be **independent** if and only if

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

That is, the probability we give A is not affected by learning the probability of B

QUESTION: If  $P(A|B) = P(A)$  can we show :

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



# Independence

Two boolean random variables A and B are said to be **independent** if and only if

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

That is, the probability we give A is not affected by learning the probability of B

QUESTION: If  $P(A)=P(A|B)$  can we show  $P(B)=P(B|A)$  ?

$$P(A \wedge B) = P(A|B) P(B) = P(A) P(B)$$

$$P(B) = \cancel{P(A \wedge B)} / \cancel{P(B)} = P(B|A)$$

$$P(A \wedge B) / P(A) \leftarrow \text{correction}$$

# Independence

If A and B are independent boolean RVs then:

- $P(A \mid B) = P(A)$  *(by definition)*
- $P(A \wedge B) = P(A \mid B) P(B) = P(A) P(B)$
- $P(B \mid A) = P(B)$

## QUESTION:

If A and B are independent boolean variables, then are their complements  $\sim A$  and  $\sim B$  also independent? That is, can we prove the following?

- $P(A \mid B) = P(A)$  if and only if  $P(\sim A \mid \sim B) = P(\sim A)$

# Independence

If A and B are independent boolean RVs, can we show their complements are independent?

Have:  $P(A)P(B) = P(A \& B)$

$$\begin{aligned} P(A \text{ or } B) &= P(A) + P(B) - P(A \& B) \\ &= P(A) + P(B) - P(A)P(B) \end{aligned}$$

$$\begin{aligned} P(\sim A \& \sim B) &= 1 - P(A \text{ or } B) \\ &= 1 - P(A) - P(B) + P(A)P(B) \\ &= (1 - P(A))(1 - P(B)) \\ &= P(\sim A)P(\sim B) \end{aligned}$$

# Multivalued Independence

- For multivalued RVs A and B, A is independent of B iff

$$\forall u, v : P(A = u \mid B = v) = P(A = u)$$

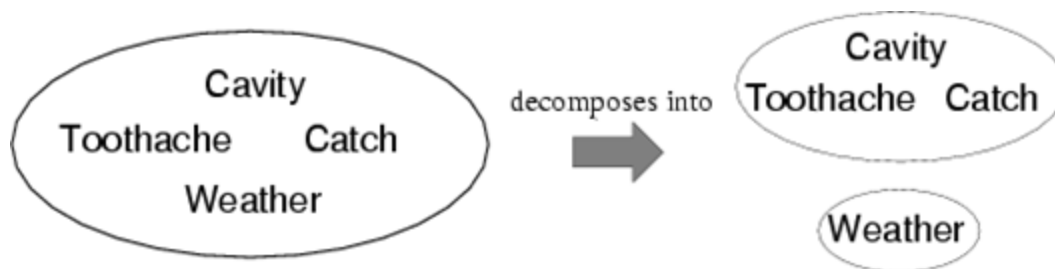
- From which we can show, for example:

$$\forall u, v : P(A = u \wedge B = v) = P(A = u)P(B = v)$$

$$\forall u, v : P(B = v \mid A = u) = P(B = v)$$

# Independence

- So, suppose our domain knowledge allows us to make certain ***independence assumptions*** on our random variables:



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

- 16 entries reduced to 10, **Why?**
- For  $n$  independent biased coins,  $O(2^n) \rightarrow O(n)$
- Absolute independence powerful but rare...
  - Dentistry is a large field with hundreds of variables, none of which are really independent of each other. **What to do?**

# Conditional independence

- "If I have a cavity, the probability that the probe catches in it doesn't depend on whether I have a toothache."

- We can then say *Catch* is conditionally independent of *Toothache* given *Cavity*:

$$(1) \mathbf{P}(\textit{catch} \mid \textit{toothache}, \textit{cavity}) = \mathbf{P}(\textit{catch} \mid \textit{cavity})$$

- The same independence holds if I haven't got a cavity:

$$(2) \mathbf{P}(\textit{catch} \mid \textit{toothache}, \neg \textit{cavity}) = \mathbf{P}(\textit{catch} \mid \neg \textit{cavity})$$

- Equivalent statements:

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

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

# Conditional independence

---

For boolean random variables,

- A is conditionally independent of B given C iff:

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

$$P(A|\sim B,C) = P(A|C)$$

For multivalued random variables,

- A is conditionally independent of B given C iff:

$$\forall u,v,w : P(A = u | B = v \wedge C = w) = P(A = u | C = w)$$

# Inference with Conditional Probabilities

---

- S = stiff neck, M = meningitis
- $P(S|M) = 0.8$ ,  $P(S) = 0.2$ ,  $P(M) = 0.0001$
- Suppose you wake up with a stiff neck - since 80% of the time, meningitis is associated with a stiff neck, you probably have meningitis and should rush to the hospital!!
- Is this correct reasoning?



# Inference with Conditional Probabilities

---

- S = stiff neck, M = meningitis
- $P(S | M) = 0.8$ ,  $P(S) = 0.2$ ,  $P(M) = 0.0001$
- $$\begin{aligned} P(M | S) &= P(M \text{ \& } S) / P(S) \\ &= P(S | M)P(M) / P(S) \\ &= (0.00008) / 0.2 \\ &= 0.0004 \end{aligned}$$
- The risk is higher, but still **very** small!

# Bayes' Theorem



## Bayes' rule:

$$P(A \mid B) = P(B \mid A) P(A) / P(B)$$

- In distribution form

$$P(Y \mid X) = P(X \mid Y) P(Y) / P(X) = \alpha P(X \mid Y) P(Y)$$

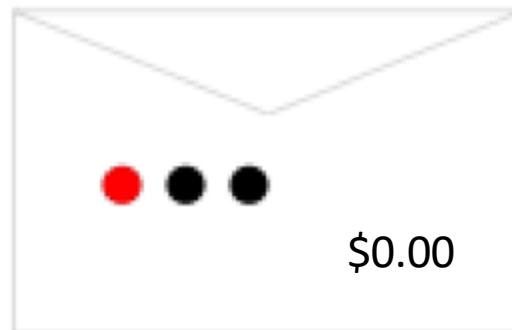
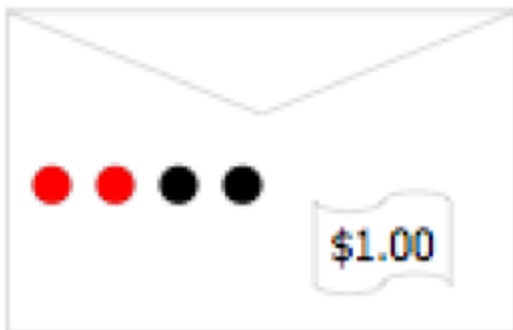
- Useful for assessing **diagnostic** probability from **causal** probability:

$$P(\text{Cause} \mid \text{Effect}) = P(\text{Effect} \mid \text{Cause}) P(\text{Cause}) / P(\text{Effect})$$

Bayes, Thomas (1783) An essay towards solving a problem in the doctrine of chances. *Philosophical Transactions of the Royal Society of London*, **53:370-418**.

# Bayes' Rule and Gambling

- Suppose there are two sealed envelopes, one (“Win”) with \$1, 2 red beads, and 2 black beads; the other with no money, 1 red bead, and 2 black beads.



- I draw an envelope at random and offer to sell it to you. How much should you be willing to pay?

# Bayes' Rule and Gambling

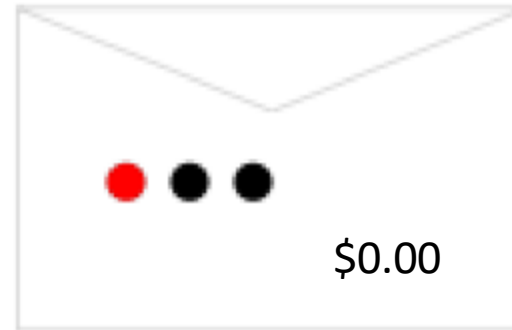
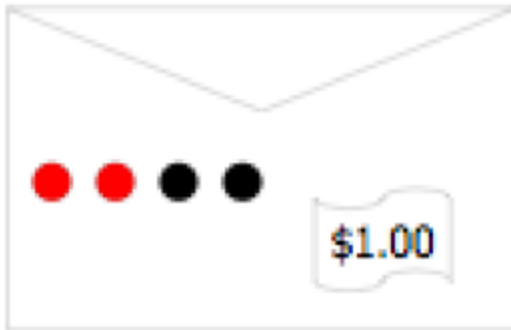
- I draw an envelope at random and offer to sell it to you. How much should you be willing to pay?



- Now, you are allowed to see one (randomly drawn) bead from the selected envelope:
  - If it is black, how much should you be willing to pay?
  - If it is red, how much should you be willing to pay?

# Bayes' Rule and Gambling

- If the bead is black...



- $P(\text{Black}) = (1/2 * 1/2 + 2/3 * 1/2) = 7/12$
- $P(\text{Win} | \text{Black}) = P(\text{Black} | \text{Win})P(\text{Win}) / P(\text{Black})$   
 $= (1/2 * 1/2) / (7/12)$   
 $= 3/7$

# Bayes' Rule and Text

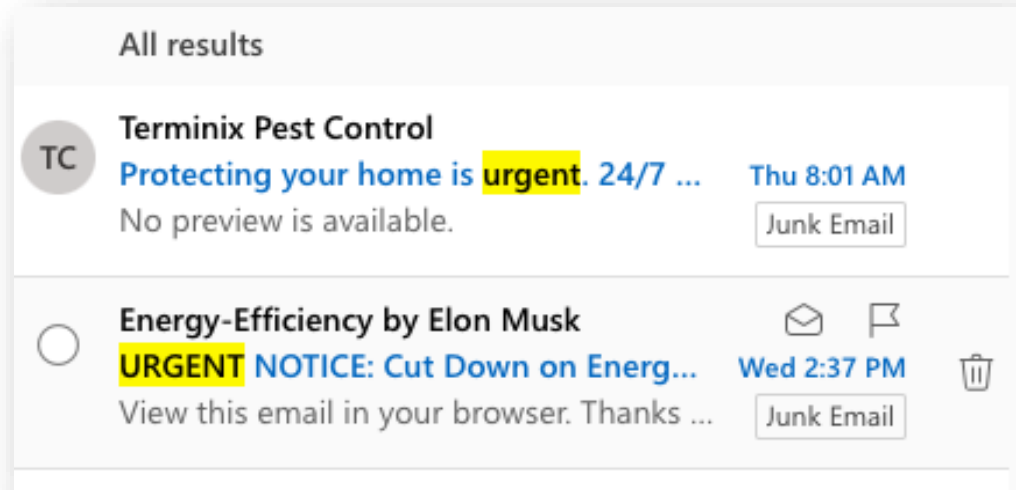
U: "Urgent" in text

S: Email is Spam

$$P(S) = 60\%$$

$$P(U|S) = 10\%$$

$$P(U|\sim S) = 1\%$$



What should we do if we receive an "Urgent" email?

# Bayes' Rule and Text

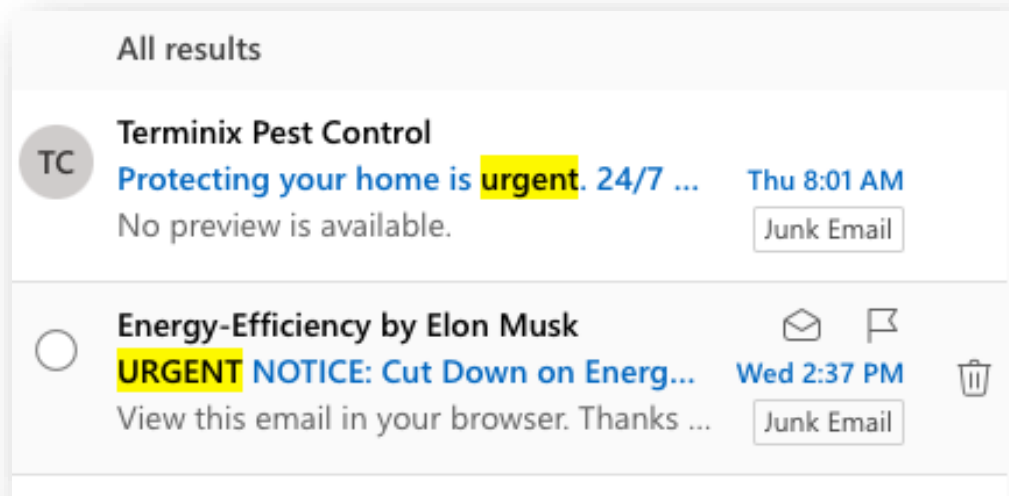
U: "Urgent" in text

S: Email is Spam

$$P(S) = 60\%$$

$$P(U|S) = 10\%$$

$$P(U|\sim S) = 1\%$$



What should we do if we receive an "Urgent" email?

$$P(U) = P(U|S)P(S) + P(U|\sim S)P(\sim S) = 0.06 + 0.004 = 0.064$$

$$\begin{aligned} P(S|U) &= P(U|S)P(S) / P(U) \\ &= 0.06 / 0.064 = \mathbf{0.9375} \end{aligned}$$

# Bayes' Rule and Text

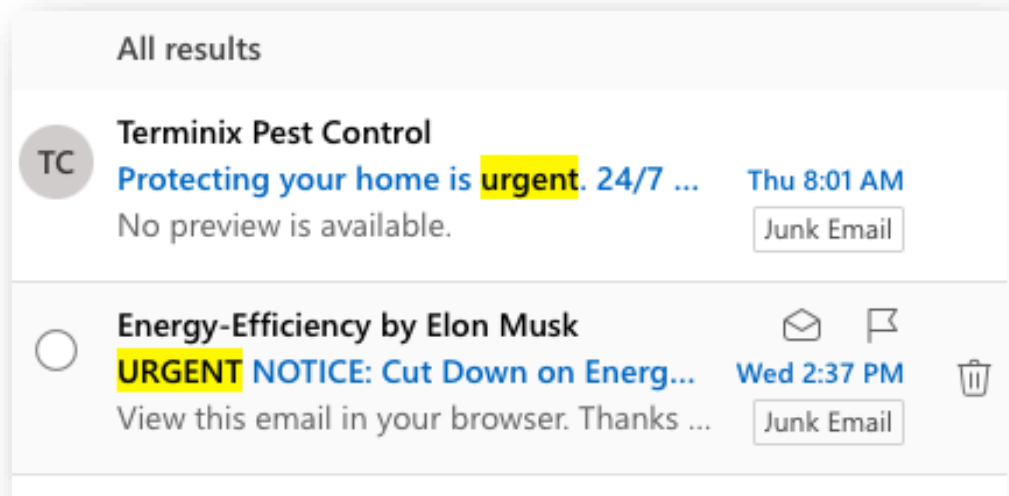
U: "Urgent" in text

S: Email is Spam

$P(S) = 60\%$  ← **Prior Prob**

$P(U|S) = 10\%$

$P(U|\sim S) = 1\%$



What should we do if we receive an "Urgent" email?

$$P(U) = P(U|S)P(S) + P(U|\sim S)P(\sim S) = 0.06 + 0.004 = 0.064$$

$$P(S|U) = P(U|S)P(S) / P(U)$$

$$= 0.06 / 0.064 = \mathbf{0.9375} \quad \leftarrow \mathbf{\text{Posterior Prob}}$$



---

# LINEAR ALGEBRA REVIEW

# Scalars, Vectors, Matrices and Tensors

- Scalars are the numbers we know and love.
- Vectors are arrays of numbers – elements of  $\mathbb{R}^n$
- They are typically written in a column (column vector)

$$\begin{aligned}a &= 1 \\b &= e \\c &= -0.3\end{aligned}$$

$$\vec{v} = \begin{bmatrix} a \\ b \\ c \end{bmatrix}$$

# Scalars, Vectors, Matrices and Tensors

- **Matrices** are sets of numbers organized into rows and columns (2-dimensional)
- Each row has the same dimension, and each column has the same dimension
- An  $M \times N$  matrix has  $M$  rows and  $N$  columns
- A vector is an  $N \times 1$  matrix
- **Tensors** are like matrices, but in higher dimensions

$$A = \begin{bmatrix} a & b \\ c & d \end{bmatrix}$$

rows

columns

# Vectors: Dot Product

---

$$a \cdot b = a^T b = \begin{bmatrix} a_1 & a_2 & a_3 \end{bmatrix} \begin{bmatrix} b_1 \\ b_2 \\ b_3 \end{bmatrix} = a_1 b_1 + a_2 b_2 + a_3 b_3$$

Think of the dot product as a matrix **multiplication**

$$\|a\|^2 = a^T a = a_1^2 + a_2^2 + a_3^2$$

The **magnitude** is the square root of the dot product of a vector with itself

$$a \cdot b = \|a\| \|b\| \cos(\theta)$$

The dot product is also related to the angle between the two vectors

# Norms

A norm is a way of measuring the magnitude of a vector

Specifically, a norm must satisfy

$$f(x) = 0 \Rightarrow x = 0$$

$$f(x + y) \leq f(x) + f(y)$$

$$f(\alpha x) = |\alpha|f(x)$$

$$L_1 \text{ Norm: } \|x\|_1 = \sum_i |x_i|$$

$$L_2 \text{ Norm: } \|x\|_2 = \sqrt{\sum_i (x_i)^2}$$

$$L_0 \text{ Norm: } \|x\|_0 = \sum_i 1 - \delta_{(x_i)(0)}$$

$$L_\infty \text{ Norm: } \|x\|_\infty = \max_i |x_i|$$

# Matrix Operations

- Addition, Subtraction, Multiplication

$$\begin{bmatrix} a & b \\ c & d \end{bmatrix} + \begin{bmatrix} e & f \\ g & h \end{bmatrix} = \begin{bmatrix} a+e & b+f \\ c+g & d+h \end{bmatrix}$$

**Just add elements**

$$\begin{bmatrix} a & b \\ c & d \end{bmatrix} - \begin{bmatrix} e & f \\ g & h \end{bmatrix} = \begin{bmatrix} a-e & b-f \\ c-g & d-h \end{bmatrix}$$

**Just subtract elements**

$$\begin{bmatrix} a & b \\ c & d \end{bmatrix} \begin{bmatrix} e & f \\ g & h \end{bmatrix} = \begin{bmatrix} ae+bg & af+bh \\ ce+dg & cf+dh \end{bmatrix}$$

**Multiply each  
row by each  
column**

# Multiplication

- Is  $AB = BA$ ? Maybe, but maybe not!

$$\begin{bmatrix} a & b \\ c & d \end{bmatrix} \begin{bmatrix} e & f \\ g & h \end{bmatrix} = \begin{bmatrix} ae + bg & \dots \\ \dots & \dots \end{bmatrix} \quad \begin{bmatrix} e & f \\ g & h \end{bmatrix} \begin{bmatrix} a & b \\ c & d \end{bmatrix} = \begin{bmatrix} ea + fc & \dots \\ \dots & \dots \end{bmatrix}$$

- Heads up: multiplication is NOT ***commutative***!

# Transpose of a Matrix

- Swap rows and columns
- The transpose of a column vector is a row vector, and vice-versa

$$A = \begin{bmatrix} a & b & c \\ d & e & f \end{bmatrix}$$

$$A^T = \begin{bmatrix} a & d \\ b & e \\ c & f \end{bmatrix}$$

$$\vec{v} = \begin{bmatrix} 1 \\ 0 \\ -1 \end{bmatrix}$$

$$\vec{v}^T = [1 \quad 0 \quad -1]$$



# Inverse of a Matrix

- Identity matrix:  
 **$AI = A$**
- Some matrices have an inverse, such that:  
 **$AA^{-1} = I$**
- Inversion is tricky:  
 **$(ABC)^{-1} = C^{-1}B^{-1}A^{-1}$**   
Derived from non-commutativity property

$$I = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

# Inverse of a Matrix

$$\begin{bmatrix} a & b & c & 1 & 0 & 0 \\ d & e & f & 0 & 1 & 0 \\ g & h & i & 0 & 0 & 1 \end{bmatrix}$$

1. Append the identity matrix to A
2. Subtract multiples of the other rows from the first row to reduce the diagonal element to 1
3. Transform the identity matrix as you go
4. When the original matrix is the identity, the identity has become the inverse!

# Orthogonality

---

- (Non-zero) vectors are **orthogonal** if their dot product is zero (geometrically, perpendicular)

$$\text{Orthogonal}(\vec{x}, \vec{y}) \stackrel{\text{def}}{=} \vec{x}^T \vec{y} = 0$$

- **Orthonormal**: orthogonal with unit norm
- An **orthogonal matrix** is one with mutually *orthonormal* rows and columns
- For an orthogonal matrix A:

$$A^{-1} = A^T$$

# Other concepts

---

- Determinant (of a matrix)
- Trace (of a matrix)
- Eigendecomposition (of a matrix)
- Pseudoinverse (of a matrix)

The matrix Cookbook:

<https://www.math.uwaterloo.ca/~hwolkowi/matrixcookbook.pdf>