# Reasoning Under Uncertainty

Uncertainty material is covered in chapters 13 and 14 of R&N and Chapter 8 of https://artint.info.

Chapter 13 gives some basic background on probability from the point of view of A.I.

Chapter 14 talks about Bayesian Networks, exact reasoning in Bayes Nets as well as approximate reasoning, which will be main topics for us.

*Note: Slides in this section draw **from** Faheim Bacchus, Craig Boutillier, Andrew Moore, Sheila McIlraith ...*

# Review: Probability Distributions over Finite Sets

A probability is a function defined over a set of atomic events U.

U represents the universe of all possible events.

# Review: Probability over Finite Sets

Given **U** (a universe of events), a probability function is a function defined over subsets of **U** that maps each subset onto the real numbers and that satisfies the Axioms of Probability. These are:

**1. P(U) = 1**

**2. P(A) $\in$ [0,1]**

**3. P({}) = 0**

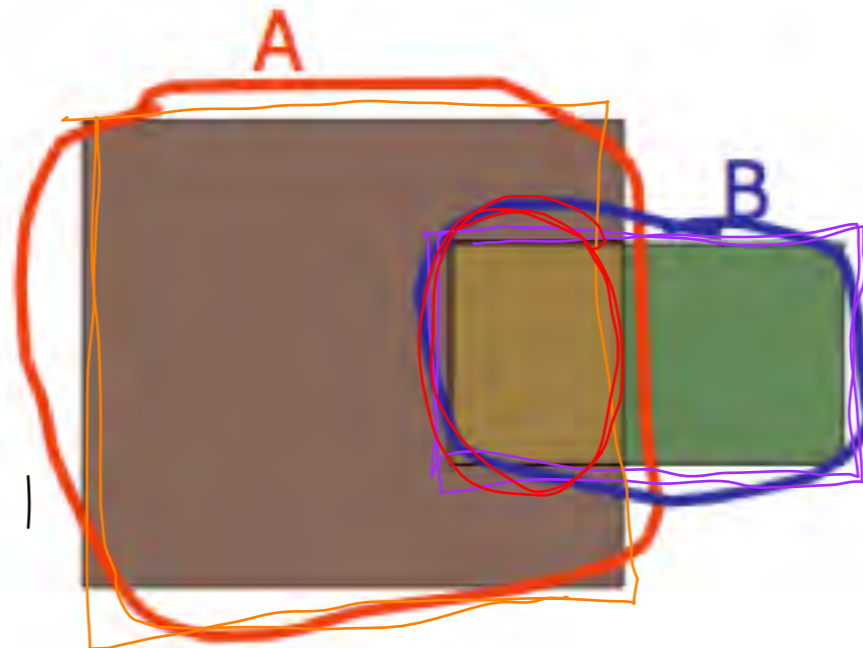**4. P(A $\cup$ B) = P(A) + P(B) – P(A $\cap$ B)**

***NB: if A $\cap$ B = {} then P(A $\cup$ B) = P(A) + P(B)***

# Review: Probability over Finite Sets

$$P(A \cup B) = P(A) + P(B) - P(A \cap B)$$

In a presidential election, there are four candidates. Call them A, B, C, and D. Based on polls, we estimate that A has a 20 percent chance of winning the election, while B has a 40 percent chance of winning. What is the probability that A or B win the election?

Vote!
http://etc.ch/ekXi

$$P(A \text{ or } B) = P(A) + P(B) - P(A \cap B)$$

$20\%$    $40\%$

# Notation: Properties and Sets

We often write

A∨B: to represent the set of events with either property A or B, i.e. the set A∪B

A∧B: to represent the set of events both property A and B, i.e. the set A∩B

¬A: to represent the set of events that do not have property A: the set U-A (i.e., the complement of A w.r.t. the universe of events U)

# Review: Probability over Feature Vectors

*CS student = {Gender, eyesight, height, eye color}*

As we move forward, ee will model sets of events in our universe as vectors of feature values.

Like CSPs, we have

1. a set of variables $V_1, V_2, \ldots, V_n$

2. a finite domain of values for each variable, $\text{Dom}[V_1], \text{Dom}[V_2], \ldots, \text{Dom}[V_n]$.

The universe of events U is the set of all vectors of values for the variables

$$\langle d_1, d_2, \ldots, d_n \rangle : d_i \in \text{Dom}[V_i]$$

When we write P(A=a, B=b), we will mean the probability that variable A has been assigned value 'a' **and** variable B has been assigned value 'b'. Note that here, sets of events are induced by a given value assignment. So, P(A=a) represents a set of events in which A holds the value 'a'.

# Review: Probability over Feature Vectors

Example:

$$P(\{V_1 = 1\}) = \sum_{x_2 \in Dom[V_2]}, \sum_{x_3 \in Dom[V_3]} P(\{V_1 = 1, V_2 = x_2, V_3 = x_3\}).$$

0.1

| | | |
|---|---|---|
| (V1 = 1, V2 = 1, V3 = 1) | (V1 = 2, V2 = 1, V3 = 1) | (V1 = 3, V2 = 1, V3 = 1) |
| (V1 = 1, V2 = 1, V3 = 2) | (V1 = 2, V2 = 1, V3 = 2) | (V1 = 3, V2 = 1, V3 = 2) |
| (V1 = 1, V2 = 1, V3 = 3) | (V1 = 2, V2 = 1, V3 = 3) | (V1 = 3, V2 = 1, V3 = 3) |
| (V1 = 1, V2 = 2, V3 = 1) | (V1 = 2, V2 = 2, V3 = 1) | (V1 = 3, V2 = 2, V3 = 1) |
| (V1 = 1, V2 = 2, V3 = 2) | (V1 = 2, V2 = 2, V3 = 2) | (V1 = 3, V2 = 2, V3 = 2) |
| (V1 = 1, V2 = 2, V3 = 3) | (V1 = 2, V2 = 2, V3 = 3) | (V1 = 3, V2 = 2, V3 = 3) |
| (V1 = 1, V2 = 3, V3 = 1) | (V1 = 2, V2 = 3, V3 = 1) | (V1 = 3, V2 = 3, V3 = 1) |
| (V1 = 1, V2 = 3, V3 = 2) | (V1 = 2, V2 = 3, V3 = 2) | (V1 = 3, V2 = 3, V3 = 2) |
| (V1 = 1, V2 = 3, V3 = 3) | (V1 = 2, V2 = 3, V3 = 3) | (V1 = 3, V2 = 3, V3 = 3) |

$U$ →

0.02

$$P(U) = 1$$

# Review: Probability over Feature Vectors

Example:

$$P(\{V_1 = 1, V_3 = 2\}) = \sum_{x_2 \in \text{Dom}[V_2]} P(\{V_1 = 1, V_2 = x_2, V_3 = 2\}).$$

(V1 = 1, V2 = 1, V3 = 1)     (V1 = 2, V2 = 1, V3 = 1)     (V1 = 3, V2 = 1, V3 = 1)
(V1 = 1, V2 = 1, V3 = 2)     (V1 = 2, V2 = 1, V3 = 2)     (V1 = 3, V2 = 1, V3 = 2)
(V1 = 1, V2 = 1, V3 = 3)     (V1 = 2, V2 = 1, V3 = 3)     (V1 = 3, V2 = 1, V3 = 3)
(V1 = 1, V2 = 2, V3 = 1)     (V1 = 2, V2 = 2, V3 = 1)     (V1 = 3, V2 = 2, V3 = 1)
(V1 = 1, V2 = 2, V3 = 2)     (V1 = 2, V2 = 2, V3 = 2)     (V1 = 3, V2 = 2, V3 = 2)
(V1 = 1, V2 = 2, V3 = 3)     (V1 = 2, V2 = 2, V3 = 3)     (V1 = 3, V2 = 2, V3 = 3)
(V1 = 1, V2 = 3, V3 = 1)     (V1 = 2, V2 = 3, V3 = 1)     (V1 = 3, V2 = 3, V3 = 1)
(V1 = 1, V2 = 3, V3 = 2)     (V1 = 2, V2 = 3, V3 = 2)     (V1 = 3, V2 = 3, V3 = 2)
(V1 = 1, V2 = 3, V3 = 3)     (V1 = 2, V2 = 3, V3 = 3)     (V1 = 3, V2 = 3, V3 = 3)

In these examples we are "summing out" some variables, which is also known as "marginalizing" our distribution

In a presidential election, there are four candidates. Call them A, B, C, and D.  Polling data and some features of the candidates are in the table below.  What's P(Winner Age > 65)? What's P(Winner Gender = Male)?

| Candidate | Gender | Region | Age | P(win) |
|-----------|--------|--------|-----|--------|
| A | Female | Midwestern | Over 65 | 0.2 |
| B | Male | Midwestern | Over 65 | 0.4 |
| C | Female | Eastern | Over 65 | 0.3 |
| D | Male | Western | Under 65 | 0.1 |

Vote!
http://etc.ch/ekXi

$$\sum_{id} \sum_{gender} \sum_{region} P(id, gender, region, age > 65)$$

# Review: Probability over Feature Vectors

Problem:

There is an exponential number of atomic probabilities to specify.

Requires summing up an exponential number of items.

To evaluate the probability of sets containing a particular subset of variable assignments we can do much better. Improvements come from the use of:

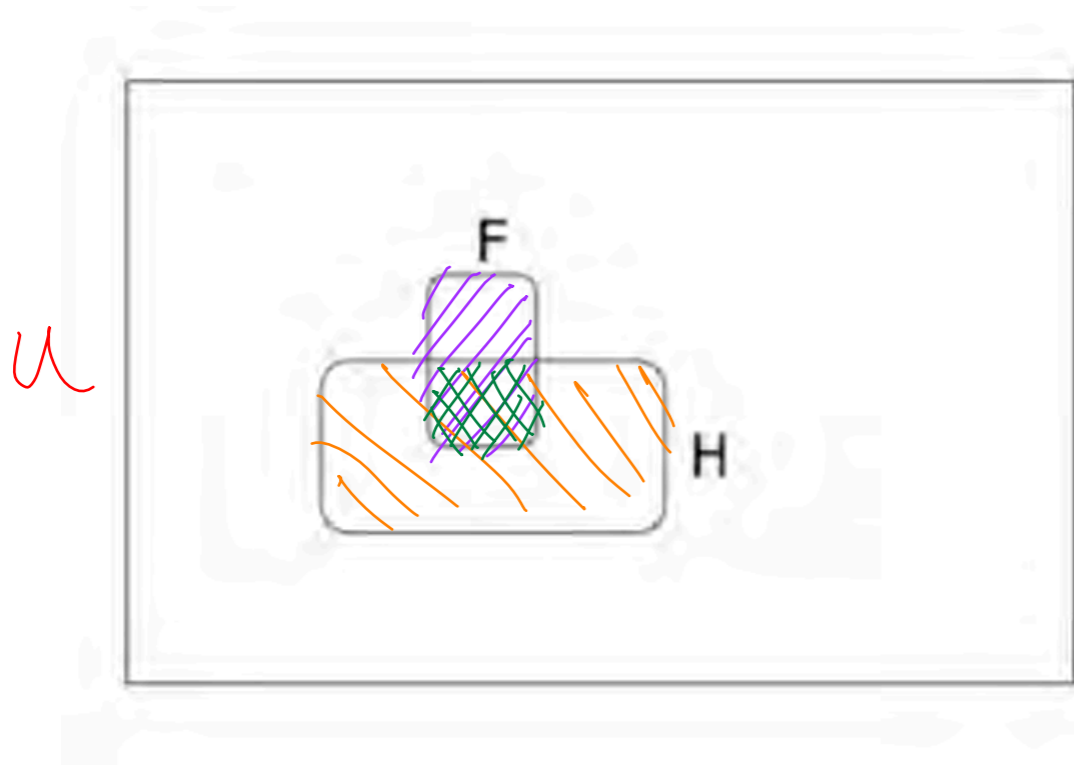1. **probabilistic independence, especially conditional independence.**

2. **approximation techniques, many of which depend on distributions structured by independence.**

# What is a Conditional Probability?

P(Headache=true|Flu=true) represents the fraction of flu-infected worlds in which you have a headache.



Green area
Purple area

= # worlds with flu and headache/#worlds with flu

= area of flu and headache/area of flu

= P(Headache=true,Flu=true)/P(Flu=true)

$$\frac{P(H=t \wedge F=t)}{P(F=t)}$$

# Axioms of Conditional Probability

A conditional probability is a also probability function, but now over a *subset* of events in the universe instead of over the entire universe. Similar axioms hold:

**P(A|A) = 1**

**P(B|A) $\in$ [0,1]**

**P(C $\cup$ B|A) = P(C|A) + P(B|A) – P(C $\cap$ B|A)**

# Review: Independence

**Probability density** is a measure of likelihood. Assume you pick an element at random from U. Density (i.e. the value of P(B) is a measure as to how likely is it to also be in set B.

It could be that the density (i.e. likelihood) of B given (or conditioned on) A is **identical** to its density (or likelihood) in U.

Alternately, the density of B given A could be very **different** that its density (or likelihood) in U.

In the first case we say that B is **independent** of A. While in the second case B is **dependent** on A.

# Review: Independence

A and B are **independent** properties:

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

A and B are **dependent**:

$$P(B|A) \neq P(B)$$

# Computational Impact

We will soon see in more detail how independence allows us to speed up computations related to inference. But the fundamental insight is that

If A and B are independent properties then

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

Proof:

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

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

# Computational Impact

We will soon see in more detail how independence allows us to speed up computations related to inference. But the fundamental insight is that

If A and B are independent properties then

**P(A∧B) = P(B) * P(A)**

Proof:

P(B|A) = P(B)                    (def'n of independence)
P(A∧B)/P(A) = P(B)
P(A∧B) = P(B) * P(A)

# What is Conditional Independence?

If $P(V_1|V_2=b, V_3=c) = P(V_1|V_2=b)$, we have not gained any additional information about $V_1$ from knowing $V_3=c$.

In this case we say that $V_1$ is conditionally independent of $V_3$ *given* $V_2$.

That is, once we know $V_2$, additionally knowing $V_3$ is irrelevant (it will give us no more information as to the value of $V_1$).

Note we could have $P(V_1|V_3=c) \neq P(V_1)$. But once we learn $V_2=b$, the value of $V_3$ becomes irrelevant.

# Computational Impact

Similar results hold for conditional independence. If B and C are conditionally independent given A, then

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

Proof:

$P(A) \cdot \dfrac{P(B, C, A)}{P(A)} =$

$P(C|A)$

$P(B|A)$

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

$\dfrac{P(B|A)}{P(C|A)} = \dfrac{P(BA)}{P(A)}$

# Computational Impact

Similar results hold for conditional independence. If B and C are conditionally independent given A, then

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

Proof:

$P(B|C \wedge A) = P(B|A)$ (def'n of conditional independence)

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

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

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

# Computational Impact

As with independence, conditional independence allows us to break up our computation onto distinct parts

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

It also allows us to ignore certain pieces of information during computations

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

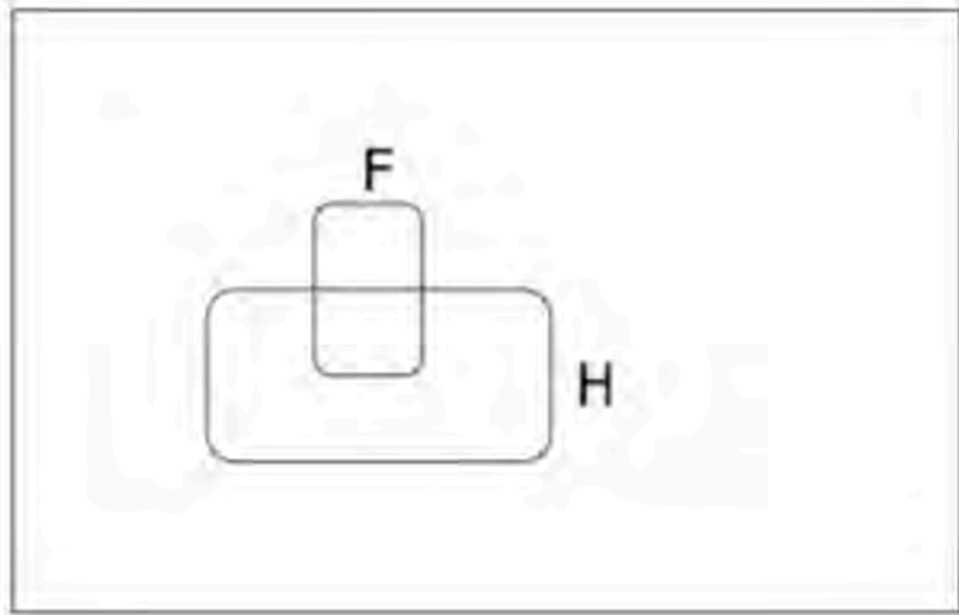# Review: Chain Rule

$$P(A_1 \wedge A_2 \wedge \ldots \wedge A_n) =$$

$$P(A_1 | A_2 \wedge \ldots \wedge A_n) * P(A_2 | A_3 \wedge \ldots \wedge A_n)$$

$$* \ldots * P(A_{n-1} | A_n) * P(A_n)$$

Proof:

$$\frac{P(A_1 \wedge A_2 \wedge \ldots \wedge A_n)}{P(A_2 \wedge \ldots A_n)} \cdot \frac{P(A_2 \wedge \ldots \wedge A_n)}{P(A_2 \wedge \ldots A_n)} \cdots \frac{P(A_{n-1} \wedge A_n)}{P(A_n)} \cdot P(A_n)$$

$$= P(A_n | A_{n-1} \wedge \ldots \wedge A_1) \cdot P(A_{n-1} | A_{n-2} \wedge \ldots A_1) \ldots \wedge P(A_1)$$

# Review: Chain Rule

$P(A_1 \wedge A_2 \wedge \ldots \wedge A_n) =$
  $P(A_1 | A_2 \wedge \ldots \wedge A_n) * P(A_2 | A_3 \wedge \ldots \wedge A_n)$
  $* \ldots * P(A_{n-1} | A_n) * P(A_n)$

Proof:

$\quad P(A_1 | A_2 \wedge \ldots \wedge A_n) * P(A_2 | A_3 \wedge \ldots \wedge A_n)$
$\; * \ldots * P(A_{n-1} | A_n)$
$= P(A_1 \wedge A_2 \wedge \ldots \wedge A_n) / P(A_2 \wedge \ldots \wedge A_n) *$
$\quad P(A_2 \wedge \ldots \wedge A_n) / P(A_3 \wedge \ldots \wedge A_n) * \ldots *$
$\quad P(A_{n-1} \wedge A_n) / P(A_n) * P(A_n)$

# Back to Flu World



P(Headache=true) = 1/10

P(Flu=true) = 1/40

P(Headache=true|Flu=true) = 1/2

Headaches are rare and having flu is rarer. But, given flu, there is a 50/50 chance you have a headache.

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

→ What is P(Flu=true|Headache=true)?

Vote!
http://etc.ch/ekXi

$$\frac{\frac{1}{2} \cdot \frac{1}{40}}{\frac{1}{10}}$$

VIEW POLL RESULTS HERE https://directpoll.com/r?XDbzPBd3ixYqg8Ave5sE9MkTtkd5vAYDdCsIT7a9h

# Vote!
## http://etc.ch/ekXi

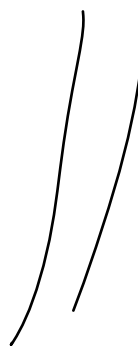$$\underset{\text{posterior}}{P(F\mid H)} = \frac{\overset{\text{prior}}{P(H\mid F)\,P(F)}}{P(H)}$$
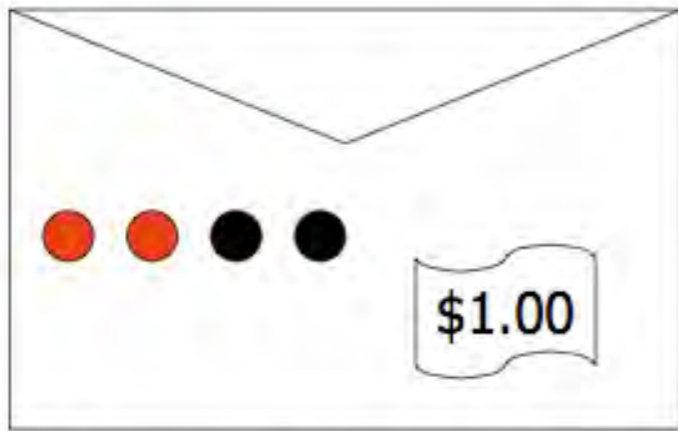
Bayes rule
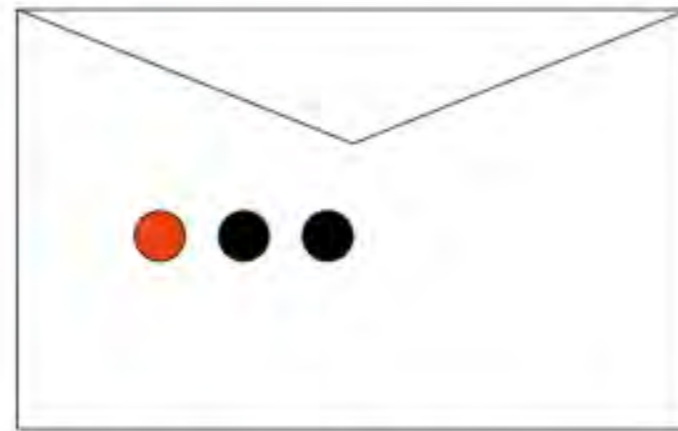
# What we just did

We Derived Bayes' Rule.

$$P(Y|X) = P(X|Y)P(Y)/P(X)$$

$$
\begin{aligned}
P(Y|X) &= P(Y \wedge X)/P(X) \\
&= P(Y \wedge X)/P(X) * P(Y)/P(Y) \\
&= P(Y \wedge X)/P(Y) * P(Y)/P(X) \\
&= P(X|Y)P(Y)/P(X)
\end{aligned}
$$

# Using Bayes Rule to gamble



The "Win" envelope has a dollar and four beads in it

The "Lose" envelope has three beads and no money

Trivial question: Someone picks an envelope and random and asks you to bet as to whether or not it holds a dollar. What are your odds?
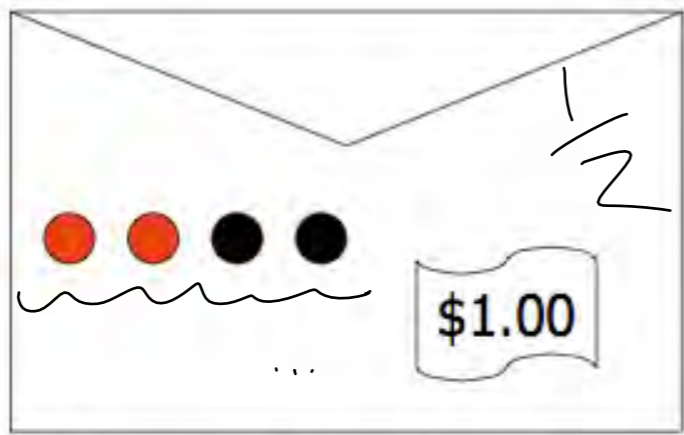
Vote!
http://etc.ch/ekXi

# Using Bayes Rule to gamble

$P(E=\$ | B=black)?$

$* \; P(B=black | E=\$) \dfrac{P(\$)}{P(B=black)}$

$P(B=black)$

$P(B=bl | E=\$)$

$\dfrac{1}{2} \cdot \dfrac{1}{2}$
_____

$\dfrac{1}{2} \cdot \dfrac{1}{2} + \dfrac{2}{3} \cdot \dfrac{1}{2}$
$\dfrac{1}{2}$ $\dfrac{1}{2}$



$\dfrac{1}{2}$

$\dfrac{2}{3}$

$= \dfrac{1}{2} \cdot \dfrac{1}{2}$

The "Win" envelope has a dollar and four beads in it

The "Lose" envelope has three beads and no money

$P(B=bl | E=\$) P(E=\$) +$

Not trivial question: Someone lets you take a bead out of the envelope before you bet. If it is black, what are your odds? If it is red, what are your odds?

$P(B=Bl | E=\overline{\$})$

$* P(E=\overline{\$})$

$\bigotimes \quad \sum_{E} P(B=bl, E) = P(B=bl)$

# Vote!
# http://etc.ch/ekXi

[VIEW POLL RESULTS HERE https://directpoll.com/r?XDbzPBd3ixYqg8Ave5sE9MkTtkd5vAYDdCsIT7a9h](https://directpoll.com/r?XDbzPBd3ixYqg8Ave5sE9MkTtkd5vAYDdCsIT7a9h)
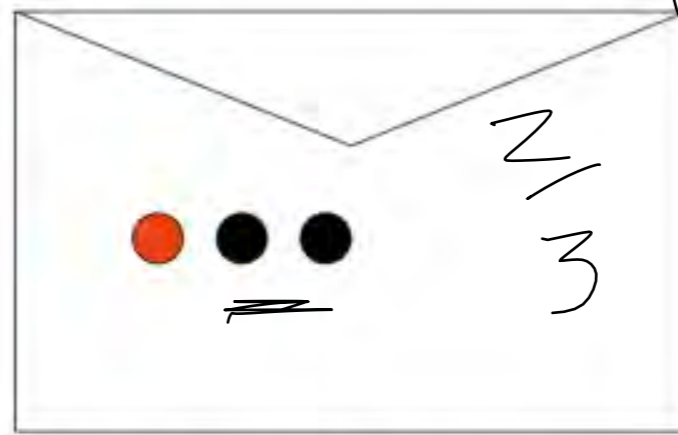
# Review: Normalizing

To **normalize** a vector of k numbers or a column in our table, e.g., <3, 4, 2.5, 1, 10, 21.5> we must sum them and divide each number by the sum:

3 + 4 + 2.5 +1 +10 + 21.5 = 42

Normalized vector:
= <3/42, 4/42, 2.5/42, 1/42, 10/42, 21.5/42>
= <0.071, 0.095, 0.060, 0.024, 0.238, 0.512>

After normalizing the vector of numbers sums to 1

It therefore can be used to specify a probability distribution.

In a presidential election, there are four candidates. Call them A, B, C, and D.  Polling data and some features of the candidates are in the table below.  Given the winning candidate is female, what's the probability that they are Midwestern?

| Candidate | Gender | Region | Age | P(win) |
| --- | --- | --- | --- | --- |
| A | Female | Midwestern | Over 65 | 0.2 |
| B | Male | Midwestern | Over 65 | 0.4 |
| C | Female | Eastern | Under 65 | 0.3 |
| D | Male | Western | Under 65 | 0.1 |

$$P(\text{Midwest} \mid \female)$$

$$\frac{0.2}{0.5}$$

Vote!
http://etc.ch/ekXi

$$P(\text{Midwest} + \female) \mid P(\female)$$

$$0.2/0.5$$