

## Project Pitches – This Thursday in Lab

- 1 minute
- Worth 1% of your mark
- Convince your classmates it's a good project idea
- Include
  - what is the task?
  - why is it an important problem?
  - how will you solve it? is there data already available?

## Probability continued

### Interpreting the axioms

- $0 \leq P(A) \leq 1$
- $P(\text{True}) = 1$
- $P(\text{False}) = 0$
- $P(A \text{ or } B) = P(A) + P(B) - P(A \text{ and } B)$

Monotonicity: if A is a subset of B, then  $P(A) \leq P(B)$

Proof:

- A subset of B  $\rightarrow B = A + C$  for  $C = B - A$
- A and C are disjoint  $\rightarrow P(B) = P(A \text{ or } C) = P(A) + P(C)$
- $P(C) \geq 0$
- So  $P(B) \geq P(A)$

### Interpreting the axioms

- $0 \leq P(A) \leq 1$
- $P(\text{True}) = 1$
- $P(\text{False}) = 0$
- $P(A \text{ or } B) = P(A) + P(B) - P(A \text{ and } B)$

Theorem:  $P(\sim A) = 1 - P(A)$

Proof:

- $P(A \text{ or } \sim A) = P(\text{True}) = 1$
- A and  $\sim A$  are disjoint  $\rightarrow P(A \text{ or } \sim A) = P(A) + P(\sim A)$
- $\rightarrow P(A) + P(\sim A) = 1$
- ....then solve for  $P(\sim A)$

## Multivalued Discrete Random Variables

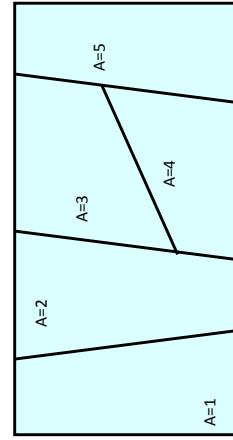
- Suppose A can take on more than 2 values
- A is a random variable with arity k if it can take on exactly one value out of  $\{v_1, v_2, \dots, v_k\}$
- Example:  $A = \{1, 2, 3, \dots, 20\}$ ; good for 20-sided dice games
- Notation: let's write the event A HasValueOf v as " $A = v$ "
- Thus...

$$P(A = v_i \wedge A = v_j) = 0 \text{ if } i \neq j$$

$$P(A = v_1 \vee A = v_2 \vee \dots \vee A = v_k)$$

## Elementary Probability in Pictures

$$\sum_{j=1}^k P(A = v_j) = 1 \quad (\text{law of total probability})$$



## Definition of Conditional Probability

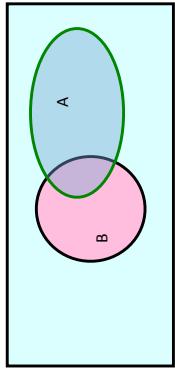
$$P(A \wedge B) \rightarrow P(A|B) = \frac{P(A \wedge B)}{P(B)}$$

We say "probability of A given b"

## Definition of Conditional Probability

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

Foundation for Bayes' Rule!



Corollary: The Chain Rule

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

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

## Independent Events

- Definition: two events A and B are *independent* if  $P(A \text{ and } B) = P(A) * P(B)$
- Intuition: knowing A tells us nothing about the value of B (and vice versa)
- From chain rule  

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

$$\rightarrow P(A|B) = P(A)$$
- You frequently need to assume the independence of something to solve a learning problem.

## Continuous Random Variables

- The discrete case: sum over all values of A is 1

- The continuous case: infinitely many values for A and the integral is 1  

$$\int_{-\infty}^{\infty} f_P(x) dx = 1$$

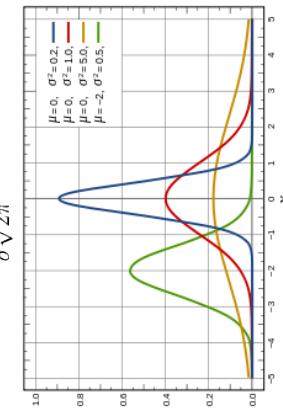
$f(x)$  is a probability density function (pdf)

1. $0 \leq P(A) \leq 1$ 2. $P(\text{True}) = 1$ 3. $P(A \text{ or } B) = P(A) + P(B) - P(A \text{ and } B)$	$\forall x, f_P(x) \geq 0$ also...
---	---------------------------------------

## Continuous Random Variables

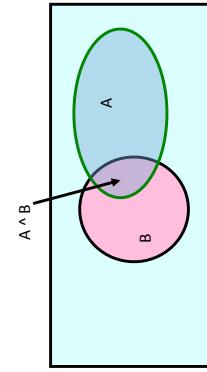
Gaussian probability density with parameters  
 - mean  $\mu$   
 - standard deviation  $\sigma$

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$



## Bayes Rule

- let's write two expressions for  $P(A \wedge B)$



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



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

we call  $P(A)$  the "prior"  
and  $P(A|B)$  the "posterior"

### Other Forms of Bayes Rule

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

Bayes, Thomas (1763) An essay towards  
solving a problem in the doctrine of chances.  
*Philosophical Transactions of the Royal Society  
of London*, 53:370-418

...by no means merely a curious speculation in the doctrine of chances, but  
necessary to be solved in order to a sure foundation for all our reasonings  
concerning past facts, and what is likely to be hereafter.... necessary to be  
considered by any that would give a clear account of the strength of *analogical* or  
*inductive reasoning*...

### Applying Bayes Rule

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

$A =$  you have the flu,  $B =$  you just coughed

Assume:

$$P(A) = 0.05$$

Also assume the following information is known to you

$$P(B|A) = 0.80$$

$$P(B|\neg A) = 0.2$$

what is  $P(\text{flu} | \text{cough}) = P(A|B)$ ?

16

### The Joint Distribution

Example: Boolean variables A,  
B, C

Recipe for making a joint distribution of M  
variables:

1. Make a truth table listing all combinations of values of your variables (if there are M Boolean variables then the table will have  $2^M$  rows).

### The Joint Distribution

Example: Boolean variables A,  
B, C

Recipe for making a joint distribution of M  
variables:

A	B	C
0	0	0
0	0	1
0	1	0
0	1	1
1	0	0
1	0	1
1	1	0
1	1	1

## The Joint Distribution

## The Joint Distribution

Example: Boolean variables A, B, C

A	B	C	Prob
0	0	0	0.30
0	0	1	0.05
0	1	0	0.10
0	1	1	0.05
1	0	0	0.05
1	0	1	0.10
1	1	0	0.25
1	1	1	0.10

Recipe for making a joint distribution of M variables:

1. Make a truth table listing all combinations of values of your variables (if there are M Boolean variables then the table will have  $2^M$  rows).
2. For each combination of values, say how probable it is.

3. If you subscribe to the axioms of probability, those numbers must sum to 1.

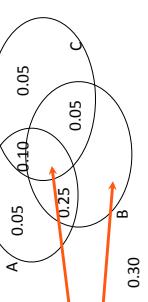
What goes here?

Example: Boolean variables A, B, C

A	B	C	Prob
0	0	0	0
0	0	1	0.05
0	1	0	0.10
0	1	1	0.05
1	0	0	0.05
1	0	1	0.10
1	1	0	0.25
1	1	1	0.10

Recipe for making a joint distribution of M variables:

1. Make a truth table listing all combinations of values of your variables (if there are M Boolean variables then the table will have  $2^M$  rows).
2. For each combination of values, say how probable it is.
3. If you subscribe to the axioms of probability, those numbers must sum to 1.



What goes here?

## Joint Probability Distribution

gender	hours_worked	wealth
Female	v0:40-5+	poor 0.253122
	rich 0.0245895	rich 0.0245895
Male	v0:40-5+	poor 0.0421768
	rich 0.0116293	rich 0.0116293

## Using the Joint Distribution

Once you have the joint distribution, you can ask for the probability of any logical expression involving your attribute

gender	hours_worked	wealth
Female	v0:40-5+	poor 0.253122
	rich 0.0245895	rich 0.0245895
Male	v0:40-5+	poor 0.0421768
	rich 0.0116293	rich 0.0116293

$$P(E) = \sum_{\text{rows matching } E} P(\text{row})$$

$$P(\text{Poor}) = 0.7604$$

## Inference with the Joint

gender	hours_worked	wealth
Female	v0:40-5+	poor 0.253122
	rich 0.0245895	rich 0.0245895
Male	v0:40-5+	poor 0.0421768
	rich 0.0116293	rich 0.0116293

## Maximum Likelihood Estimation (MLE)

Rich vs Poor

$$P(E_1 | E_2) = \frac{P(E_1 \wedge E_2)}{P(E_2)} = \frac{\sum_{\text{rows matching } E_1 \text{ and } E_2} P(\text{row})}{\sum_{\text{rows matching } E_2} P(\text{row})}$$

What is the probability of a person being rich, given you know nothing else about that person?



## Why 3/5?

We assume that the wealth of the people in our dataset  $D$  is independently distributed

$$\theta = \text{Probability of being rich} = P(\text{rich})$$

$$? = \text{Probability of being poor} = P(\text{poor})$$

$$D = \{r, p, r, r, p\}$$

$$\alpha_r = \# \text{ rich}$$

$$\alpha_p = \# \text{ poor}$$

$$P(D) = P(r \text{ and } p \text{ and } r \text{ and } r \text{ and } p)$$

## A little math

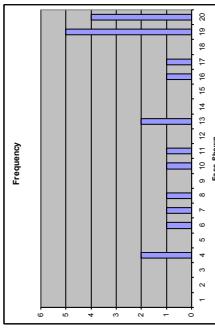
$$\underset{\theta}{\operatorname{argmax}} P(D) = (1 - \theta)^{\alpha_F} * \theta^{\alpha_H}$$

## That's Maximum Likelihood Estimation (MLE)

It's not always the best solution...

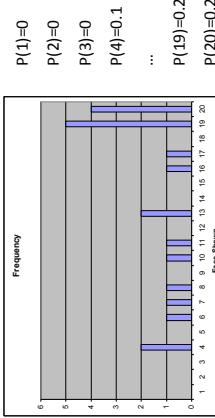
## Issues with MLE estimate

I bought a loaded d20 on eBay...but it didn't come with any specs. How can I find out how it behaves?



## Issues with MLE estimate

I bought a loaded d20 on eBay...but it didn't come with any specs. How can I find out how it behaves?

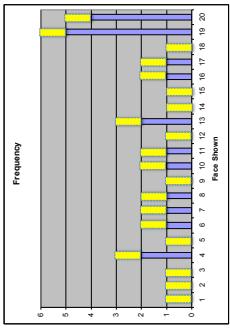


1. Collect some data (20 rolls)
2. Estimate  $P(i) = \text{CountOf(rolls of } i\text{)} / \text{CountOf(any roll)}$

But: Do I really think it's *impossible* to roll a 1,2 or 3?

## A better solution

I bought a loaded d20 on eBay...but it didn't come with any specs.  
How can I find out how it behaves?



0. Imagine some data (20 rolls, each i shows up 1x)

1. Collect some data (20 rolls)
2. Estimate P(i)

What if we know that poor people are much more common than rich people?



We have a belief about  $\theta$

- $P(\theta | D) = P(D | \theta) * P(\theta) / P(D)$

- $\propto P(D | \theta) * P(\theta)$

Now we can incorporate our belief about  $\theta$

This is a MAP (Maximum A Posteriori) Estimate

## Conjugate Prior

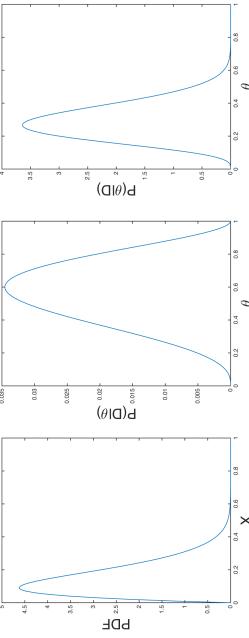
- Our likelihood so far has been based on a Bernoulli distribution.
- Beta is a conjugate prior to Bernoulli
  - This means their pdfs play nice together
  - $P(D | \theta) * P(\theta)$  will be easy to deal with
  - Called the posterior likelihood

## Beta/Binomial Distributions

- Beta
  - $P(\theta) \propto (1 - \theta)^{\beta_p - 1} * \theta^{\beta_r - 1}$
  - proportional to (missing a constant)
- Binomial

$$P(\theta) = (1 - \theta)^{\alpha_p} * \theta^{\alpha_r}$$

## More math



$$P(\theta | D) P(\theta) \propto [(1 - \theta)^{\alpha_p} * \theta^{\alpha_r}] * [(1 - \theta)^{\beta_p - 1} * \theta^{\beta_r - 1}]$$

Exercises:

1. find the optimal  $\theta$
2. What if  $\beta_p = 1$  and  $\beta_r = 1$ ?

## Estimating Parameters

- Maximum Likelihood Estimate (MLE): choose  $\theta$  that maximizes probability of observed data

$$\hat{\theta} = \arg \max_{\theta} P(\mathcal{D} | \theta)$$

- Maximum a Posteriori (MAP) estimate: choose  $\theta$  that is most probable given **prior probability and the data**

$$\hat{\theta} = \arg \max_{\theta} P(\theta | \mathcal{D}) \\ = \arg \max_{\theta} \frac{P(\mathcal{D} | \theta) P(\theta)}{P(\mathcal{D})}$$

A wonderful tutorial:  
<http://www.cs.ubc.ca/~murphyk/Teaching/CS340-Fall06/reading/bernoulli.pdf>