

Lecture 1: Introduction

Julia Hockenmaier

juliahmr@illinois.edu
3324 Siebel Center

<http://www.cs.uiuc.edu/class/sp10/cs598jhm>

Why should you take this class?

Many recent NLP papers use **Bayesian methods**.

If you want to do research in NLP, you need to understand these papers!

You may wonder:

- What are Bayesian methods?

(or: what's Gibbs sampling, variational inference, LDA, HDP, conjugate priors?)

- Why are they used?

- How can you use them in your own research?

You may also wonder:

- What's the difference between generative and discriminative models, or generative and discriminative estimation?

(and why should I care?)

Why should you take this class?

CS546 has been canceled!
(and you need to fill your schedule)

What are we going to do in this class?

Graphical models are the foundation of much of machine learning used in NLP (and elsewhere).

Graphical models enable us to work with very complex probabilistic models, because they have efficient **inference** and **parameter estimation** methods.

You've (hopefully) all seen graphical models before, but it won't harm to go over the basics again and delve a bit deeper.

What are we going to do in this class?

Foundations:

Relevant parts of probability theory
Graphical models
Parameter estimation
Inference: variational methods and sampling

Applications:

Topic models (LDA etc.)
Finite-state methods: e.g. HMMs and extensions
Grammar-based methods

What are we going to do in this class?

Lectures:

Background material

Paper presentations:

Current research

Readings:

Background material and current papers

Homework:

More reading (and some writing)

Research project:

Applying what we're learning (and again some writing)

Grades etc.

Your grade will consist of:

- 50% **research project**
- 30% **class presentations**
- 10% (mini) **literature reviews**
- 10% **class participation**

The research project

The goal: write a research paper.

We may submit good papers to conferences/workshops.
You can work alone or in pairs.

The project can be related to your thesis research, but it needs to be clear that what you are doing for this class is something different from what you would have done anyway.

Grading criteria:

- technical soundness
- originality
- quality of writeup

You may have to give a couple of brief presentations in class about your project.

Class presentations

You will be expected to present **two conference papers** (or **one longer journal paper**) in class.

Grading criteria:

- Clarity of presentation
- Can you interpret the paper?

- We will upload your slides to the course website.

Literature reviews

There will be **five homeworks**, in which you will be asked to write a 1-2 page minireview of a conference paper.

What was done? Why? How?

What would be a good follow-up paper?

Class website; email

Slides, pointers to reading material etc. will be posted at:

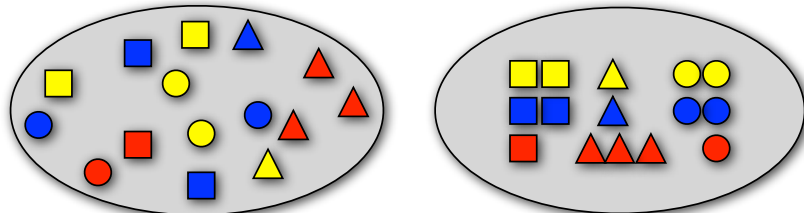
<http://www.cs.uiuc.edu/class/sp10/cs598jhm>

We'll also have a mailing list.

Basic Probability Theory

Sampling with replacement

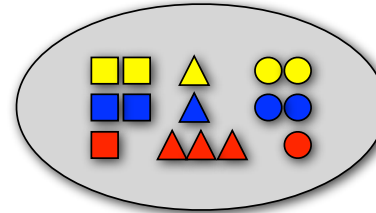
Pick a random shape, then put it back in the bag.



$$\begin{aligned}
 P(\square) &= 2/15 & P(\square) &= 1/15 & P(\square \text{ or } \triangle) &= 2/15 \\
 P(\text{blue}) &= 5/15 & P(\text{red}) &= 5/15 & P(\triangle | \text{red}) &= 3/5 \\
 P(\text{blue} | \square) &= 2/5 & P(\square) &= 5/15 & &
 \end{aligned}$$

Sampling with replacement

Pick a random shape, then put it back in the bag.
What **sequence of shapes** will you draw?



$$\begin{aligned}
 P(\circ \triangle \triangle \square) &= 1/15 \times 1/15 \times 1/15 \times 2/15 \\
 &= 2/50625
 \end{aligned}$$

$$\begin{aligned}
 P(\triangle \circ \circ \triangle) &= 3/15 \times 2/15 \times 2/15 \times 3/15 \\
 &= 36/50625
 \end{aligned}$$

$$\begin{aligned}
 P(\square) &= 2/15 & P(\triangle) &= 1/15 & P(\circ) &= 2/15 \\
 P(\square) &= 2/15 & P(\triangle) &= 1/15 & P(\circ) &= 2/15 \\
 P(\square) &= 1/15 & P(\triangle) &= 3/15 & P(\circ) &= 1/15
 \end{aligned}$$

Sampling with replacement

Alice was beginning to get very tired of sitting by her sister on the bank, and of having nothing to do: once or twice she had peeped into the book her sister was reading, but it had no pictures or conversations in it, 'and what is the use of a book,' thought Alice 'without pictures or conversation?'

$$\begin{aligned}
 P(\text{of}) &= 3/66 & P(\text{her}) &= 2/66 \\
 P(\text{Alice}) &= 2/66 & P(\text{sister}) &= 2/66 \\
 P(\text{was}) &= 2/66 & P(,) &= 4/66 \\
 P(\text{to}) &= 2/66 & P(') &= 4/66
 \end{aligned}$$

Sampling with replacement

beginning by, very Alice but was and? reading no tired of to into sitting sister the, bank, and thought of without her nothing: having conversations Alice once do or on she it get the book her had peeped was conversation it pictures or sister in, 'what is the use had twice of a book' 'pictures or' to

$$\begin{aligned}
 P(\text{of}) &= 3/66 & P(\text{her}) &= 2/66 \\
 P(\text{Alice}) &= 2/66 & P(\text{sister}) &= 2/66 \\
 P(\text{was}) &= 2/66 & P(,) &= 4/66 \\
 P(\text{to}) &= 2/66 & P(') &= 4/66
 \end{aligned}$$

In our model, $P(\text{English}) = P(\text{word salad})$

Conditioning on the previous word

Alice was beginning to get very tired of sitting by her sister on the bank, and of having nothing to do: once or twice she had peeped into the book her sister was reading, but it had no pictures or conversations in it, 'and what is the use of a book,' thought Alice 'without pictures or conversation?'

$$\begin{aligned} P(w_{i+1} = \text{of} \mid w_i = \text{tired}) &= 1 & P(w_{i+1} = \text{bank} \mid w_i = \text{the}) &= 1/3 \\ P(w_{i+1} = \text{of} \mid w_i = \text{use}) &= 1 & P(w_{i+1} = \text{book} \mid w_i = \text{the}) &= 1/3 \\ P(w_{i+1} = \text{sister} \mid w_i = \text{her}) &= 1 & P(w_{i+1} = \text{use} \mid w_i = \text{the}) &= 1/3 \\ P(w_{i+1} = \text{beginning} \mid w_i = \text{was}) &= 1/2 \\ P(w_{i+1} = \text{reading} \mid w_i = \text{was}) &= 1/2 \end{aligned}$$

Probability theory: terminology

Trial: picking a shape, predicting a word

Sample space Ω : the set of all possible outcomes (all shapes; all words in *Alice in Wonderland*)

Event $\omega \subseteq \Omega$: an actual outcome (a subset of Ω) (predicting 'the', picking a triangle)

The probability of events

Kolmogorov axioms:

$$\begin{aligned} 0 &\leq P(\omega \subseteq \Omega) \leq 1 \\ P(\emptyset) &= 0 \text{ and } P(\Omega) = 1 \\ \sum_{\omega_i \subseteq \Omega} P(\omega_i) &= 1 \quad \text{if } \forall j \neq i : \omega_i \cap \omega_j = \emptyset \\ &\quad \text{and } \bigcup_i \omega_i = \Omega \end{aligned}$$

The ω_i form a partition of Ω

Random variables

A random variable X is a function which maps every element in the sample space to some value.

Discrete random variables:

- heads or tails,
- age (in years), day of the week,
- number of letters/vowels/e's in a word,
- part of speech of a word
- the word itself

Continuous random variables:

- size, height, weight,

Discrete probability distributions: Throwing a coin

Bernoulli distribution:

Probability of success (=head,yes) in single yes/no trial

- The probability of *head* is p .
- The probability of *tail* is $1-p$.

Binomial distribution:

Prob. of the number of heads in a sequence of yes/no trials

The probability of getting exactly k heads in n independent yes/no trials is:

$$P(k \text{ heads, } n - k \text{ tails}) = \binom{n}{k} p^k (1 - p)^{n-k}$$

Discrete probability distributions: Rolling a die

Categorical distribution:

Probability of getting one of N outcomes in a single trial

- The probability of category/outcome c_i is p_i ($\sum p_i = 1$)

Multinomial distribution:

Probability of observing each possible outcome c_i exactly X_i times in a sequence of n yes/no trials

$$P(X_1 = x_1, \dots, X_N = x_N) = \frac{n!}{x_1! \dots x_N!} p_1^{x_1} \dots p_N^{x_N} \quad \text{if } \sum_{i=1}^N x_i = n$$

Expectation

The **expectation** $E(X)$ of a *discrete* random variable X with probability $p(x=X)$ is

$$E(X) = \sum_x p(x)x$$

$E(X)$ is also called the **expected value, mean or average**.

The expectation $E(X)$ of a *continuous* random variable X with probability density function $p(x=X)$ is

$$E(X) = \int_{-\infty}^{+\infty} p(x)x dx$$

X =number of vowels/word

Words:

lab, think, eat, come, mail, student, book, lecture, coffee, facebook

$$p(X=1)=0.2 \quad p(X=2)=0.5$$

$$p(X=3)=0.2 \quad p(X=4)=0.1$$

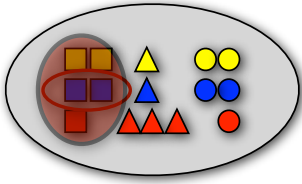
$$\begin{aligned} E(X) &= 1 \times 0.2 + 2 \times 0.5 + 3 \times 0.2 + 4 \times 0.1 \\ &= 2.2 \end{aligned}$$

The expected value of a random variable may not correspond to any possible value!

Joint and Conditional Probability

The **conditional probability of X given Y** , $P(X|Y)$, is defined in terms of the probability of Y , $P(Y)$, and the **joint probability of X and Y** , $P(X,Y)$:

$$P(X|Y) = \frac{P(X,Y)}{P(Y)}$$



$$P(\text{blue} | \blacksquare) = 2/5$$

The chain rule

Conversely, the joint probability $P(X,Y)$ can also be expressed in terms of the conditional probability $P(X|Y)$

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

This leads to the so-called **chain rule**:

$$\begin{aligned} P(X_1, X_2, \dots, X_n) &= P(X_1)P(X_2|X_1)P(X_3|X_2, X_1)\dots P(X_n|X_1, \dots, X_{n-1}) \\ &= P(X_1) \prod_{i=2}^n P(X_i|X_1 \dots X_{i-1}) \end{aligned}$$

Independence

Two random variables X and Y are independent if

$$P(X, Y) = P(X)P(Y)$$

If X and Y are independent, then $P(X|Y) = P(X)$:

$$\begin{aligned} P(X|Y) &= \frac{P(X, Y)}{P(Y)} \\ &= \frac{P(X)P(Y)}{P(Y)} \quad (X, Y \text{ independent}) \\ &= P(X) \end{aligned}$$

Bayes' theorem

If the events A_1, \dots, A_k form a partition of the space Ω such that $P(A_j) > 0$ for $j=1, \dots, k$, and B is any event such that $P(B) > 0$, then, for $i=1, \dots, k$, then:

$$P(A_i|B) = \frac{P(A_i)P(B|A_i)}{\sum_{j=1}^k P(A_j)P(B|A_j)}$$

Example

We have some widgets,

- 20% were produced by machine M1
- 30% were produced by machine M2
- 50% were produced by machine M3

Unfortunately,

- 1% of the items produced by M1 are defective
- 2% of the items produced by M2 are defective
- 3% of the items produced by M3 are defective

We have one defective widget.

What's the probability that it was produced by M2?

Priors and posteriors

$$P(A_i|B) = \frac{P(A_i)P(B|A_i)}{\sum_{j=1}^k P(A_j)P(B|A_j)}$$

$P(A_i)$ is the **prior** probability of A_i

$P(A_i|B)$ is the **posterior** probability of A_i

Bayesian statistics

- We use θ to represent (the parameters of) probability distributions
- Probabilities θ represent degrees of belief
- Data D provides evidence for or against our beliefs.
We update our belief θ based on the evidence we see:

$$P(\theta|D) = \frac{P(\theta)P(D|\theta)}{\int P(\theta)P(D|\theta)d\theta}$$

- For a given set of data D , $P(D|\theta)$ is the **likelihood** of θ