# CS447: Natural Language Processing

*http://courses.engr.illinois.edu/cs447*

# Lecture 1: Introduction

## Julia Hockenmaier

*juliahmr@illinois.edu*

3324 Siebel Center

# Course Staff

Professor:
  Julia Hockenmaier   juliahmr@illinois.edu

Teaching assistants:
  Qingrong Cheng     qc16@illinois.edu
  Maghav Kumar   mkumar10@illinois.edu

# Today's lecture

Course Overview:

What is NLP?
What will you learn in this course?

Course Admin:

How will we teach this course?
How will you be assessed in this course?

# What is Natural Language Processing?

# What is Natural Language Processing *really* ?

# NLP in the news…

Facebook AI Creates Its Own Language In Creepy Preview Of Our Potential Future

Computers can now describe images using language you'd understand

## Barbie Wants to Get to Know Your Child

With the help of A.I., America's most famous doll tries to fulfill a timeless dream — convincing little girls that she's a real friend. What will happen if they believe her?

# Have chatbots developed their own language?

*https://code.facebook.com/posts/1686672014972296/deal-or-no-deal-training-ai-bots-to-negotiate/*

*Bob: "I can can I I everything else."*
*Alice: "Balls have zero to me to me to me to me to me to me to me to me to."*

No, not really…
http://www.snopes.com/facebook-ai-developed-own-language/

# How well does that Barbie do?

Barbie: *"Do you have any sisters?"*
Child:   "Yeah, I only have one."
Barbie: *"What's something nice that your sister does for you?"*
Child:   "She does nothing nice to me"
Barbie: *"Well, what is the last nice thing your sister did?"*
Child:   "She helped me with my project
               — and then she destroyed it."
Barbie:  *"Oh, yeah, tell me more!"*
Child:    "That's it, Barbie,"
Barbie:  *"Have you told your sister lately how cool she is?"*
Child:    "No. She is not cool,"
Barbie:  *"You never know, she might appreciate hearing it"*

(Thanks to Barbara Grosz for pointing out the example from the NYT)
https://www.nytimes.com/2015/09/20/magazine/barbie-wants-to-get-to-know-your-child.html

# IBM's Watson wins at Jeopardy!

# Machine Translation



Google Translate

# Dialog systems, chatbots, digital assistants

# Siri [Jurafsky & Martin, 3rd ed]



(a)                                                                    (b)

Challenge: Dialogue consists of several turns.
Systems need to be able to understand and sanswer
follow-up questions, etc.

# What is the current state of NLP?

Lots of commercial applications and interest.
  Some applications are working pretty well already,
  others not so much.

A paradigm shift (and a lot of hype) around "deep
learning" and "AI" more generally
  —Neural nets are powerful classifiers and sequence models
  —Public libraries (Tensorflow, Pytorch, etc.) and datasets
  make it easy for anybody to get a model up and running
  —"End-to-end" models put into question whether we still need
  the traditional NLP pipeline that this class is built around
  —We're still in the middle of this paradigm shift
  —But many of the fundamental problems haven't gone away

# What will you learn in this class?

# What will you learn in this class?

In this class, you will learn about some of

— the core **tasks** (as well as **data sets and evaluation metrics**) that people work on in NLP

— the fundamental **models and algorithms** that have been developed for these tasks

— the relevant **linguistic concepts and phenomena** that will be encountered in these tasks

# The focus of this class

We want to identify the structure and meaning of words, sentences, texts and conversations
N.B.: we do not deal with speech (no signal processing)

We mainly deal with language analysis/understanding, and less with language generation/production

We focus on fundamental concepts, methods, models, and algorithms, not so much on current research:
- Data (natural language): linguistic concepts and phenomena
- Representations: grammars, automata, etc.
- Neural and statistical models over these representations
- Learning & inference algorithms for these models

# What you should learn

You should be able to answer the following questions:
- What makes natural language difficult for computers?
- What are the core NLP tasks?
- What are the main modeling techniques used in NLP?

We won't be able to cover the latest research…
(this requires more time, and a much stronger background in machine learning than I am able to assume for this class)

… but I would still like you to get an understanding of:
- How well does current NLP technology work (or not)?
- What NLP software and datasets are available?
- How to read NLP research papers [4 credits section]

# Building a computer that '*understands*' text: The traditional NLP pipeline

新华社拉萨二月二日电（记者央珍） "八五"（一九九一至一九九五年）期间，西藏金融体制改革坚持与全国框架一致、体制衔接的方针，顺利完成了西藏各级人民银行的分设工作，实现信贷资金使用从粗放型经营方式向集约型经营方式转变。

# Task: Tokenization/segmentation

新华社拉萨二月二日电（记者央珍）"八五"（一九九一至一九九五年）期间，西藏金融体制改革坚持与全国框架一致、体制衔接的方针，顺利完成了西藏各级人民银行的分设工作，实现信贷资金使用从粗放型经营方式向集约型经营方式转变。

We need to split text into words and sentences.

- Languages like Chinese don't have spaces between words.
- Even in English, this cannot be done deterministically:

*There was an earthquake near D.C. You could even feel it in Philadelphia, New York, etc.*

NLP task:

What is the *most likely* segmentation/tokenization?

# Task: Part-of-speech-tagging

Open the pod door, Hal.

⬇

| Verb | Det | Noun | Noun | , | Name | . |
|------|-----|------|------|---|------|---|
| **Open** | **the** | **pod** | **door** | **,** | **Hal** | **.** |

**open:**
verb, adjective, or noun?
Verb: **open** *the door*
Adjective: *the* **open** *door*
Noun: *in the* **open**

# How do we decide?

We want to know the most likely tags $T$
for the sentence $S$

$$\underset{T}{\operatorname{argmax}} P(T|S)$$

We need to define *a statistical model* of $P(T \mid S)$, e.g.:

$$\underset{T}{\operatorname{argmax}} P(T|S) \quad = \quad \underset{T}{\operatorname{argmax}} P(T)P(S|T)$$

$$P(T) \quad =_{def} \quad \prod_i P(t_i|t_{i-1})$$

$$P(S|T) \quad =_{def} \quad \prod_i P(w_i \mid t_i)$$

We need to estimate *the parameters* of $P(T \mid S)$, e.g.:

$$P(\,t_i = V \mid t_{i-1} = N\,) = 0.3$$

# Disambiguation requires statistical models

Ambiguity is a core problem for any NLP task

Statistical models* are one of the main tools
to deal with ambiguity.

*more generally: a lot of the models (classifiers, structured prediction models)
you learn about in CS446 (Machine Learning) can be used for this purpose.
You can learn more about the connection to machine learning in CS546
(Machine learning in Natural Language).

These models need to be trained (estimated, learned)
before they can be used (tested).

We will see lots of examples in this class
(CS446 is NOT a prerequisite for CS447)

# *"I made her duck"*

What does this sentence mean?
  *"duck"*: noun or verb?
  *"make"*: *"cook X"* or *"cause X to do Y"* ?
  *"her"*: *"for her"* or *"belonging to her"* ?

Language has different kinds of ambiguity, e.g.:
## Structural ambiguity
  *"I eat sushi **with tuna**"* vs. *"I eat sushi **with chopsticks**"*
  *"I saw the man **with the telescope on the hill**"*
## Lexical (word sense) ambiguity
  *"I went to the **bank**"*:  financial institution or river bank?
## Referential ambiguity
  *"**John** saw **Jim. He** was drinking coffee."*

# *"I made her duck cassoulet"*

(Cassoulet = a French bean casserole)

The second major problem in NLP is **coverage**:
We will always encounter unfamiliar words
and constructions.

Our models need to be able to deal with this.

This means that our models need to be able
to *generalize* from what they have been trained on
to what they will be used on.

# Task: Syntactic parsing

# Observation: Structure corresponds to meaning



**Correct analysis**

VP
   NP
     PP
V   NP   P   NP
eat   sushi   with tuna

eat sushi with tuna

VP
  VP     PP
V   NP   P   NP
eat  sushi  with chopsticks

eat sushi with chopsticks

**Incorrect analysis**

VP
  VP     PP
V   NP   P   NP
eat  sushi  with tuna

eat sushi with tuna

VP
   NP
     PP
V   NP   P   NP
eat  sushi  with chopsticks

eat sushi with chopsticks

# Question: what is grammar?

**Grammar formalisms**

**(= linguists' programming languages)**

A precise way to define and describe
the structure of sentences.

**Specific grammars**

**(= linguists' programs)**

Implementations (in a particular formalism)
for a particular language (English, Chinese,....)

Overgeneration

English

Undergeneration

John Mary saw.

with tuna sushi ate I.

Did you went there?

....

Did you go there?

I want you to go there.

I ate the cake that John had made for me yesterday

John and Mary eat sushi for dinner.

.....

John saw Mary.

I ate sushi with tuna.

# NLP and automata theory

What kind of grammar/automaton
is required to analyze natural language?

What class of languages does
natural language fall into?

Chomsky (1956)'s hierarchy of formal languages
was originally developed to answer (some of)
these questions.

# Task: Semantic analysis

$$\exists x \exists y (\text{pod\_door}(x) \ \& \ \text{Hal}(y)$$
$$\& \ \text{request}(\text{open}(x, y)))$$

```
                    S
          VP
              NP
                  NOUN                    NP
                                           |
    Verb   Det  Noun Noun ,  Name .
    Open   the  pod  door ,  Hal    .
```

# Representing meaning

We need a meaning representation language.

**"Shallow" semantic analysis:** Template-filling (Information Extraction)
  Named-Entity Extraction: Organizations, Locations, Dates,...
  Event Extraction

**"Deep" semantic analysis:** (Variants of) formal logic

`∃x∃y(pod_door(x)& Hal(y) & request(open(x,y)))`

We also distinguish between
**Lexical semantics** (the meaning of words) and
**Compositional semantics** (the meaning of sentences)

# Understanding texts

More than a decade ago, Carl Lewis stood on the threshold of what was to become the greatest athletics career in history. He had just broken two of the legendary Jesse Owens' college records, but never believed he would become a corporate icon, the focus of hundreds of millions of dollars in advertising. His sport was still nominally amateur. Eighteen Olympic and World Championship gold medals and 21 world records later, Lewis has become the richest man in the history of track and field -- a multi-millionaire.

Who is Carl Lewis?
Did Carl Lewis break any world records?
(and how do you know that?)

# Multimodal NLP: mapping from language to the world

$$\exists x \exists y (\texttt{pod\_door(x)} \ \& \ \texttt{Hal(y)}$$
$$\& \ \texttt{request(open(x, y)))}$$

**System**

`request(open(`**`door2`**`, `**`SYS`**`))`

# Summary: The NLP Pipeline

An NLP system may use some or all
of the following steps:

Tokenizer/Segmenter
  to identify words and sentences
Morphological analyzer/POS-tagger
  to identify the part of speech and structure of words
Word sense disambiguation
  to identify the meaning of words
Syntactic/semantic Parser
  to obtain the structure and meaning of sentences
Coreference resolution/discourse model
  to keep track of the various entities and events mentioned

# NLP Pipeline: Assumptions

Each step in the NLP pipeline embellishes the input with **explicit information** about its linguistic structure

POS tagging: parts of speech of word,

Syntactic parsing: grammatical structure of sentence,….

Each step in the NLP pipeline requires its own explicit **("symbolic") output representation**:

POS tagging requires a POS tag set

(e.g. NN=common noun singular, NNS = common noun plural, …)

Syntactic parsing requires constituent or dependency labels

(e.g. NP = noun phrase, or nsubj = nominal subject)

These representations should capture linguistically appropriate **generalizations/abstractions**

Designing these representations requires linguistic expertise

# NLP Pipeline: Shortcomings

Each step in the pipeline relies on a **learned model** that will return the *most likely* representations

- This requires a lot of **annotated training data** for each step
- Annotation is **expensive** and sometimes **difficult** (people are not 100% accurate)
- These models are **never 100% accurate**
- Models make more mistakes if their input contains mistakes

How do we know that we have captured the "*right*" **generalizations** when designing representations?

- Some representations are **easier to predict** than others
- Some representations are **more useful** for the next steps in the pipeline than others
- But we won't know how easy/useful a representation is until we have a model that we can plug into a particular pipeline

# Course Admin

# This class consists of...

## … Lectures:

Wednesdays and Fridays, 12:30pm–1:45 pm, DCL1310

## … Office Hours:

Julia: Wednesdays and Fridays, 2:00pm–2:30pm, Siebel 3324

Qingrong: TBD, Siebel 0207

Meghav: TBD, Siebel 0207

## … Websites:

Syllabus, slides, policies, etc: *http://courses.engr.illinois.edu/cs447*

Discussions: *piazza.com/illinois/fall2019/cs447*

Grades, submitting assignments: Compass2G *http://compass2g.illinois.edu*
and Gradescope http://gradescope.com

## … Readings:

Textbook + additional readings (*http://courses.engr.illinois.edu/cs447*)

## … Assessment:

4 assignments, 2 exams (4th credit hour: project or survey)

# Lectures and office hours

Attend!

Ask questions!

# Participate!

# Reading

**Course website: (slides, reading)**

https://courses.engr.illinois.edu/cs447/fa2018/syllabus.html

**The textbook:** https://web.stanford.edu/~jurafsky/slp3/

Jurafsky and Martin, **Speech and Language Processing**
(We will mostly use chapters from the forthcoming 3rd
edition, but may also use the 2nd edition from 2008 in a few
place)

# Assessment

If you take this class for 3 hours credit:
   1/3 homework assignments
   1/3 midterm exam
   1/3 final exam

If you take this class for 4 hours credit:
   1/4 homework assignments
   1/4 midterm exam
   1/4 final exam
   1/4 literature review or project

We reserve the right to improve your grade by up to 5% depending on your class participation. If you're in between grades, but attended class and participated frequently and actively in in-class discussions etc., we will give you the higher grade.

# Homework assignments

## What?

4 assignments (mostly programming)

We use Python 3

## Why?

To make sure you can put what you've learned to practice.

## How?

You will have three weeks to complete HW1, HW2, HW3, HW4.

Grades will be based on your write-up and your code.

Submit your assignments on Gradescope.

## Late policy?

**No** late assignments will be accepted (sorry).

# Homework assignments

Schedule:
Week 2:   Friday, 09/06   HW1 out
Week 5:   Friday, 09/27   HW1 due, HW2 out
Week 8:   Friday, 10/20   HW2 due, HW3 out
Week 11: Friday, 11/08   HW3 due, HW4 out
Week 14: Friday, 12/06   HW4 due

Points per assignment:
HW1,HW2,HW3,HW4 = 10 points per assignment

# Exams

## What?

First midterm: Friday, Oct 11, in class
Second Midterm: Wednesday, Dec 11, in class
   (based on material after first midterm)

## Why?

To make sure you understand what you learned
well enough to explain and apply it.

## How?

Essay questions and problem questions
Closed-book (no cheatsheets, no electronics, etc.)
Will be based on lectures and readings

# 4th credit hour: Research Projects

## What?

You need to read and describe a few (2–3) NLP papers on a particular task, implement an NLP system for this task and describe it in a written report.

## Why?

To make sure you get a deeper knowledge of NLP by reading original papers and by building an actual system.

## When?

**Fri, Oct 4:** Proposal due (What topic? What papers will you read?)
**Fri, Nov 8:** Progress report due (Are your experiments on track?)
**Thu, Dec 12:** Final report due (Summary of papers, your system)

# 4th credit hour: Literature Survey

### What?

You need to read and describe several (5-7) NLP papers on a particular task or topic, and produce a written report that compares and critiques these approaches.

### Why?

To make sure you get a deeper knowledge of NLP by reading original papers, even if you don't build an actual system.

### When?

**Fri, Oct 4:** Proposal due (What topic? What papers will you read?)
**Fri, Nov 8:** Progress report due (Is your paper on track?)
**Thu, Dec 12:** Final report due (Summary of papers)

# Today's readings

**Today's lecture:**

Jurafsky and Martin Chapter 1 (2nd edition)

http://www.cs.colorado.edu/~martin/SLP/Updates/1.pdf