# CS 4650/7650: Natural Language Processing

Jacob Eisenstein

Lecture 1: Introduction

August 17, 2015

- Lecture: KACB 1456, Monday/Wednesday 3:05-4:30
- Office hours (Eisenstein): CCB 316, Thursday 10-11.
- My email: jacobe@gatech
- TA: TBA
- Webpage:
  `https://github.com/jacobeisenstein/gt-nlp-class/`

# Prerequisites

- **Officially**: CS 3510, Design and Analysis of Algorithms
- **Unofficially**
    - Basic linear algebra
    - Solid probability and statistics
    - Automata and formal language theory: e.g., finite-state vs context-free languages, etc
    - Ability to analyze and implement dynamic programming algorithms
    - Coding ability (Python strongly preferred)
    - Helpful, but not assumed:
        - Some familiarity with basic machine learning:
          naïve Bayes, logistic regression, perceptron
        - Some familiarity with basic ideas about linguistics

Readings should be completed before the lecture on the date assigned.
The will be drawn from:

- **My notes**
- **Linguistic fundamentals for NLP** by Emily Bender.
- **Foundations of Statistical NLP** by Manning and Schuetze
- Other online resources.

# Resources

Recommended background reading:

- Jurafsky and Martin, Second Edition
- **The NLTK Book** by Bird, Klein, Loper
- **Linguistic Structure Prediction** by Noah Smith
- **Machine Learning** by Kevin Murphy
- **Introduction to Information Retrieval** by Manning, Raghavan, & Schütze
- **Probability: The Analysis of Data, Vol 1** by Guy Lebanon
- Journals: Computational Linguistics, Journal of Machine Learning Research, Transactions of the Association of Computational Linguistics (TACL)
- Conferences: ACL, NAACL, EMNLP, EACL, NIPS, ICML, ...

Assignments

- Six assigned problem sets (48%)
- Twelve **short** homework assignments. You may skip two. (20%)
- In-class midterm exam (12%)
- Final exam (20%)

Assignments

- Six assigned problem sets (48%)
- Twelve **short** homework assignments. You may skip two. (20%)
- In-class midterm exam (12%)
- Final exam (20%)

There will be a lot of reading, a lot of coding, and a lot of math.

Through the problem sets, you will:

- Build increasingly complex and practical NLP systems.
- Test properties of these systems by performing experiments.
- Derive properties of NLP systems mathematically.
- Compete in in-class "bakeoff" competitions.
- A student: "The best parts were the projects, which encompassed the complete spectrum of NLP."

# Problem set grading

- Assignments are due at the **beginning** of lecture.
- Accepted up to 3 days late, with a penalty of 20% per day.
- A detailed collaboration policy is online here `https://github.com/jacobeisenstein/gt-nlp-class/blob/master/Grading.md`
- I take academic integrity very seriously.
  See `www.honor.gatech.edu` and the online syllabus for more details.
  If you have a question about this policy, ask!

Through the **short** homework assignments you will:

- Learn to identify linguistic phenomena by labeling real texts.
- Compare your linguistic analysis with your classmates.
- Help critique your classmates' independent projects.
- Not spend more than one hour per assignment, usually less than thirty minutes.

# Homeworks

Through the **short** homework assignments you will:

- Learn to identify linguistic phenomena by labeling real texts.
- Compare your linguistic analysis with your classmates.
- Help critique your classmates' independent projects.
- Not spend more than one hour per assignment, usually less than thirty minutes.

Grading

- Submit a PDF online by the beginning of class. If possible, please bring a paper copy to class to discuss.
- There will be twelve homeworks. You may skip two.
- Homeworks will not be accepted late.
- You must work alone.

# Midterm exam

There will be an in-class midterm on October 19.

- Barring an institute approved absence, you must take the exam in class on October 19.
- The purpose of the midterm is to test your understanding of the concepts covered in class and in the readings.
- The secondary purpose is to encourage you to review those concepts.
- The midterm will include **anything** covered in class and in the readings through October 19.

Hours per week of work: $\sim 12$

- 6-9 hours: 3
- 9-12 hours: 3
- 12-15 hours: 4
- 15-18 hours: 3

# Here's what people are saying about this course

Hours per week of work: $\sim 12$

- 6-9 hours: 3
- 9-12 hours: 3
- 12-15 hours: 4
- 15-18 hours: 3

- "I expended more effort in this course than expected."
- "This class requires much more work than the average graduate class. Would be great to emphasize that during class introduction."

# In the beginning



While he was inventing the field of AI, Alan Turing asked: "how do we know when we're done?"

While he was inventing the field of AI, Alan Turing asked: "how do we know when we're done?"

The Turing Test: Can a computer carry on a conversation so naturally that you can't distinguish it from a human?

Turing, 1950. "Computing Machinery and Intelligence." Mind (236): 433-460.

```
http:
//www.pandorabots.com/pandora/talk?botid=f5d922d97e345aa1
https://www.youtube.com/watch?v=vphmJEpLXU0
```

# The Turing Test today

```
http:
//www.pandorabots.com/pandora/talk?botid=f5d922d97e345aa1
https://www.youtube.com/watch?v=vphmJEpLXU0
```

- The best chatbots today avoid deep language understanding and focus on exhaustive string matching.
- In contrast, most of NLP is concerned with building software that understands language on a deeper level.

```
http:
//www.pandorabots.com/pandora/talk?botid=f5d922d97e345aa1
https://www.youtube.com/watch?v=vphmJEpLXU0
```

- The best chatbots today avoid deep language understanding and focus on exhaustive string matching.
- In contrast, most of NLP is concerned with building software that understands language on a deeper level. **why is this hard?**

Some real examples:

- Iraqi head seeks arms

Some real examples:

- Iraqi head seeks arms
- Teacher strikes idle kids

# Ambiguity

Some real examples:

- Iraqi head seeks arms
- Teacher strikes idle kids
- Ban on nude dancing on governor's desk

# Ambiguity

Some real examples:

- Iraqi head seeks arms
- Teacher strikes idle kids
- Ban on nude dancing on governor's desk
- Stolen painting found by tree

## Ambiguity

Some real examples:

- Iraqi head seeks arms
- Teacher strikes idle kids
- Ban on nude dancing on governor's desk
- Stolen painting found by tree

# Ambiguity

Some real examples:

- Iraqi head seeks arms
- Teacher strikes idle kids
- Ban on nude dancing on governor's desk
- Stolen painting found by tree

Sometimes ambiguity is not so obvious.

- Our company is training workers.

# Ambiguity

Some real examples:

- Iraqi head seeks arms
- Teacher strikes idle kids
- Ban on nude dancing on governor's desk
- Stolen painting found by tree

Sometimes ambiguity is not so obvious.

- Our company is training workers.

Ambiguity grows with sentence length, sometimes exponentially.

A range of disambiguating cues:

- **The lexicon** (words and their syntactic functions)
  Teacher is almost always a noun, never a verb

A range of disambiguating cues:

- **The lexicon** (words and their syntactic functions)
  Teacher is almost always a noun, never a verb
- **Syntax**
  Sentences rarely have two adjacent verbs (if they are both indicative)
  The sequence ADJECTIVE-VERB-END is unlikely in English

## How?

A range of disambiguating cues:

- **The lexicon** (words and their syntactic functions)
  Teacher is almost always a noun, never a verb

- **Syntax**
  Sentences rarely have two adjacent verbs (if they are both indicative)
  The sequence ADJECTIVE-VERB-END is unlikely in English

- **Lexical semantics**
  A ban cannot dance
  A head (the body part) rarely seeks

A range of disambiguating cues:

- **The lexicon** (words and their syntactic functions)
  Teacher is almost always a noun, never a verb

- **Syntax**
  Sentences rarely have two adjacent verbs (if they are both indicative)
  The sequence ADJECTIVE-VERB-END is unlikely in English

- **Lexical semantics**
  A ban cannot dance
  A head (the body part) rarely seeks

- "**Common sense**"
  Teachers aren't supposed to hit kids

- Writing down all of these constraints and preferences in a single coherent representation is hard.
- Many (most?) sentences won't satisfy all constraints.
  How to decide which ones can be safely ignored?

- Writing down all of these constraints and preferences in a single coherent representation is hard.
- Many (most?) sentences won't satisfy all constraints.
  How to decide which ones can be safely ignored?
- The answer is data, and probability.

NLP research may still be as ambitious as the Turing test.

NLP research may still be as ambitious as the Turing test.
But it may also be very down-to-earth…

- Finding the price of products on the web
- Analyzing reading level or authorship
- Detecting sentiment about products, stocks, or world leaders
- Extracting facts or relations from documents

# Application: document classification



Email is reliably separated into priority, regular, and spam.

# Application: content and subjectivity analysis

Realtime Coverage

**Obama pushes change on historic Myanmar visit**
The News International - 1 hour ago
YANGON: President Barack Obama urged Myanmar on Monday to hasten its "remarkable" reforms on a historic visit during which he was feted by huge crowds and met Aung San Suu Kyi at the home where she was long locked up. The trip, the first to ...

For Obama and Clinton, Their Final Tour as Partners
New York Times - 15 minutes ago

President Obama on Burma tour
The University of Hawaii Kaleo - 50 minutes ago

Obama gets warm welcome in historic trip to Myanmar
New York Daily News - 58 minutes ago

Asia trip takes Obama White House into Myanmar time warp
Free Malaysia Today - 1 hour ago

Obama Makes History With Myanmar, Cambodia...
ABC News - 1 hour ago

# Application: content and subjectivity analysis

## In Depth

**For Obama and Clinton, Their Final Tour as Partners**
New York Times - 15 minutes ago
PHNOM PENH, Cambodia — They emerged from Air Force One together, side by side, smiling at the crowd waiting on the tarmac below. Then as they headed down the stairs, she held back just a little so that she would stay a step behind him.

**In a Changing Myanmar, Vows of Support From a Visiting President**
New York Times - 2 hours ago
YANGON, Myanmar — President Obama journeyed to this storied tropical outpost of jade and jungles on Monday to "extend the hand of friendship" as a land long tormented by repression and poverty begins to throw off military rule and emerge from decades ...

**Obama meets Aung San Suu Kyi**
Irish Times - 2 hours ago
irishtimes.com - Last Updated: Monday, November 19, 2012, 10:12. Obama meets Aung San Suu Kyi. US president Barack Obama kisses Aung San Suu Kyi following joint remarks at her residence in Yangon, Burma today. Photograph: Jason Reed/Reuters ...

# Application: content and subjectivity analysis

## Opinion

**The Irish Times - Tuesday, November 20, 2012**
Irish Times - 1 hour ago
Aung San Suu Kyi's caution is understandable and justified by her own experience of false dawns. Speaking at Barack Obama's side as he visited Rangoon yesterday, the Burmese opposition leader and Nobel prize winner warned of her country's tentative ...

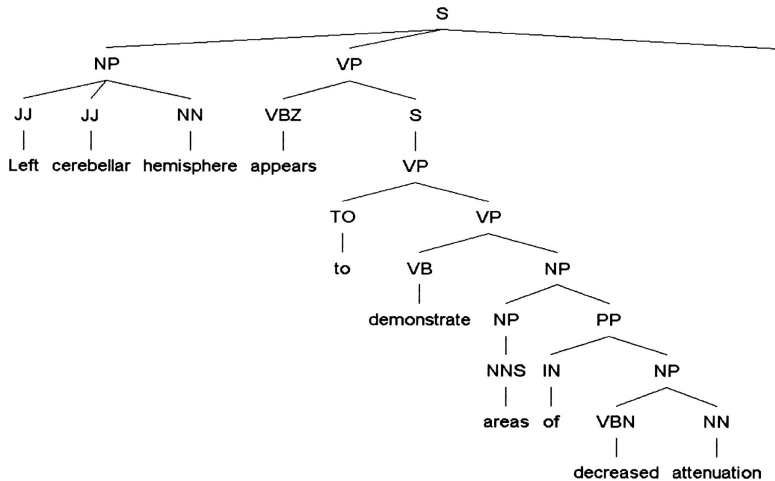**Bush's Burma Policy, Obama's Victory Lap**
Wall Street Journal - Nov 18, 2012
In one of those gems that reveal the Obama administration's penchant for taking credit for the work of others, a senior State Department official on a plane to Perth last week for a U.S.-Australia confab spoke to reporters about the president's trip to Burma ...

**President Obama Goes to Asia**
New York Times - Nov 16, 2012
President Obama leaves on Saturday for a trip to Asia that will show his commitment to having the United States engage more intensely with countries there. But it comes at an awkward time. Israel and Hamas are at war in Gaza, and efforts to end the violence ...

# Syntactic analysis today



Modern syntactic parsers get 90% accuracy on English newstext.

# Machine translation today



Le 21 décembre 2011 — *Par* clumsy

J'ai eu beaucoup de mal à trouver le point commun des albums qui m'ont hanté en 2011. Je les ai tous réécoutés, disséqués, digérés mais rien ne venait. Et puis le fil conducteur s'est dessiné. Il est devenu de plus en plus clair. De plus en plus évident : ces disques ont parlé à

Le 21 décembre 2011 — *Par* clumsy

J'ai eu beaucoup de mal à trouver le point commun des albums qui m'ont hanté en 2011. Je les ai tous réécoutés, disséqués, digérés mais rien ne venait. Et puis le fil conducteur s'est dessiné. Il est devenu de plus en plus clair. De plus en plus évident : ces disques ont parlé à

I struggled to find the common point of albums that haunted me in 2011. I've replayed all, dissected, digested, but nothing came. And the son non thread has emerged. It has become increasingly clear. More ave talked to my instinct more than s, violently attacked me, caressed me in the direction of the hair too. They've invaded the arteries and
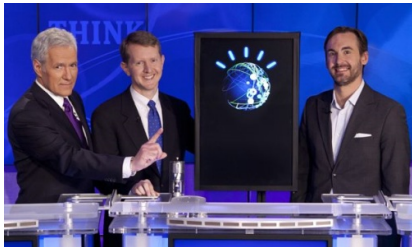
**Original French text:**       Google ⊠

J'ai eu beaucoup de mal à trouver le point commun des albums qui m'ont hanté en 2011.

⊞ Contribute a better translation

Fast, accurate, and (somewhat) fluent translation for many language pairs

Watson extracts facts from millions of documents, parses complex questions, and outperforms the best human players.

# Information extraction today



Watson extracts facts from millions of documents, parses complex questions, and outperforms the best human players.

This goes way beyond search and string match. An example question:

Wanted for general evilness, last seen at the Tower of Barad-Dur.
It's a giant eye, folks, kinda hard to miss

- All of these success stories result from applying **statistical machine learning** to large amounts of linguistic data.
- This data-driven approach will be the focus of this course.

A **corpus** is a collection of text:
often annotated in some way, but sometimes just lots of text.

The development of large corpora made
data-driven NLP possible. Some examples:

- Brown corpus:
  1M words of text with part-of-speech tags

- Penn Treebank:
  1M words of text with parse trees

- Europarl:
  1.8M aligned French-English sentence pairs

- Google n-grams:
  1.2B 5-grams and their counts
  (to think about: why is this useful?)

# Corpora

A **corpus** is a collection of text:
often annotated in some way, but sometimes just lots of text.

The development of large corpora made
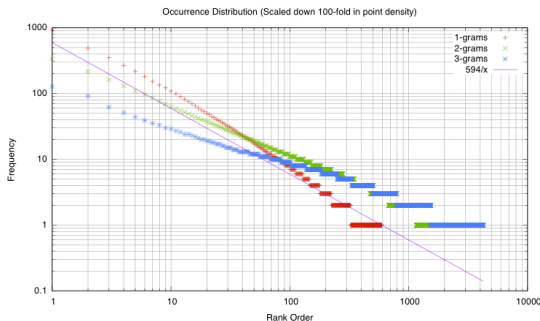data-driven NLP possible. Some examples:

- Brown corpus:
  1M words of text with part-of-speech tags

- Penn Treebank:
  1M words of text with parse trees

- Europarl:
  1.8M aligned French-English sentence pairs

- Google n-grams:
  1.2B 5-grams and their counts
  (to think about: why is this useful?)

How much data is enough?

# How much is enough?

Answer: there's no data like more data.



- There is always a "long tail" of rare but important phenomena.
- We are increasingly interested in languages with few resources, and new problems without annotations.

Natural language processing applications are typically built in a **stack**

- From "low-level" phenomena like words and morphemes...
- ... to "high-level" phenomena like semantics and discourse.

Outline of topics

- **Words**: text classification, language models, morphology
- **Sequences**: hidden Markov models, part-of-speech tagging
- **Trees**: context free grammars, parsing
- **General graphs**: semantics and discourse
- **Learning**: unsupervised and semi-supervised methods
- **Applications**: translation, information extraction, dialogue

Outline of topics

- **Words**: text classification, language models, morphology
- **Sequences**: hidden Markov models, part-of-speech tagging
- **Trees**: context free grammars, parsing
- **General graphs**: semantics and discourse
- **Learning**: unsupervised and semi-supervised methods
- **Applications**: translation, information extraction, dialogue

In each section, we will cover linguistic issues, computational representations, and statistical techniques.
You will build software that put these ideas into practice.

# Course goals

By the end of the semester, you should have learned:

- What are the range of linguistic phenomena we need to address to build useful language technology.
- How to select linguistic representations that are appropriate for the problem you want to solve.
- How to apply modern machine learning techniques to solve language processing problems.
- What are the existing resources (software and data) that can help.

CS 7650: you will also learn to read current research papers in the field.

(if there's time)

- Start **problem set 1**.
- **Homework 1**: identify ambiguous sentences in the news
- Read Chapter 1 of **Linguistic Fundamentals for NLP**, if you haven't already.
- Read Chapter 3 of my notes.
- Optional supplementary reading:
  - Section 2.1 of Foundations of Statistical NLP
  - Survey on word sense disambiguation
  - LXMLS lab guide