# CS 4650/7650: Natural Language Processing

Jacob Eisenstein

Lecture 1: Introduction

August 20, 2013

- Lecture: ES&T L1105, Tuesday/Thursday 3:05-4:30
- Office hours (Eisenstein): TSRB 228A, Wednesday 11-12
- My email: jacobe@gatech
- TA: Yangfeng Ji. Office hours: Friday 11-12, outside TSRB 228
- Webpage:
  `https://github.com/jacobeisenstein/gt-nlp-class/`

# Prerequisites

- **Officially**: CS 3600, Intro to Artificial Intelligence
- **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
  - Strong coding ability (Python strongly preferred)
  - Bonuses:
    - Some familiarity with basic machine learning: naïve Bayes, logistic regression, perceptron
    - Some familiarity with basic ideas about linguistics

# Resources

Readings are are drawn from online resources. They should be completed before the lecture on date assigned.

# Resources

Readings are are drawn from online resources. They should be completed before the lecture on date assigned.

**Places to find additional material**

- Optional textbook: Jurafsky and Martin, **Second Edition**
- Emily Bender, **Linguistic fundamentals for NLP**
- Noah Smith, **Linguistic Structure Prediction**
- Kevin Murphy, **Machine Learning**
- Manning, Raghavan, & Schütze, **Introduction to Information Retrieval**
- Guy Lebanon, **Probability: The Analysis of Data, Vol 1**
- Journals: Computational Linguistics, Journal of Machine Learning Research
- Conferences: ACL, NAACL, EMNLP, EACL, NIPS, ICML, ...

# Assignments

Assignments

- Four assigned projects (40%)
- Twelve **short** homework assignments (20%, includes class participation)
- In-class midterm exam (20%)
- Independent project (20%)

Assignments

- Four assigned projects (40%)
- Twelve **short** homework assignments (20%, includes class participation)
- In-class midterm exam (20%)
- Independent project (20%)

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

- Through the assigned projects, 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.

- Through the assigned projects, 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.
- Grading
    - Assignments are due at the **beginning** of lecture on the date indicated.
    - They'll be accepted up to 3 days late, with a penalty of 20% per day.
    - You must work alone.
    - I take academic integrity very seriously.
      See `www.honor.gatech.edu` and the online syllabus for more details.
      **If you have a question, 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.

# 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.
- Grading
  - Submit a PDF online by the beginning of class **and bring a paper copy to class**.
  - Each homework is worth 2 points: one for the writeup and one for the class discussion. If you miss the discussion, you lose the point.
  - There will be twelve homeworks. You may skip two.
  - Homeworks will not be accepted late.
  - You must work alone.

There will be an in-class midterm on October 8.

- Barring a medical emergency, you must take the exam in class on that date.
- 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 8. There will not be a practice exam or study guide.

- The independent final project will give you an opportunity to delve more deeply into an area of interest.
- You may work in teams of up to three.
- The outcome may be:
  - a piece of software
  - an experimental or theoretical analysis
  - a comprehensive review of some area of research...
- I am happy to discuss possible projects with you.
  Start thinking about this early.

- **After the midterm**: start thinking about your project if you haven't already.
- November 12: present your proposal idea in class
- December 3-5: present your project results in class
- December 6: Submit initial project report
- December 13 (est): Submit final project report

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

# In the beginning



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
http://www.youtube.com/watch?v=WnzlbyTZsQY
```

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

- 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
http://www.youtube.com/watch?v=WnzlbyTZsQY
```

- 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

# Ambiguity

Some real examples:

- Iraqi head seeks arms
- Teacher strikes idle kids

Some real examples:

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

Some real examples:

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

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 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

# 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

- "**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

# Application: content and subjectivity analysis

# 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 warne 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 senio 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 ...

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

# 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

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 ... has become increasingly clear. More ... ave talked to my instinct more than ..., 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 analysis to large amounts of linguistic data.
- This data-driven approach will be the focus of this course.

# 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?)

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

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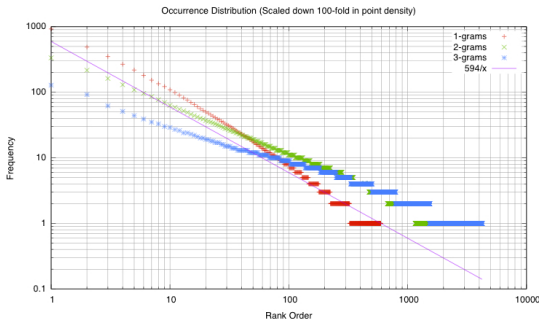data-driven NLP possible. Some examples:

- 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.
- This means we have to be smart.

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

# This course

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.

You will also learn to read current research papers in the field.

# For next time

- Homework 1: identify ambiguous sentences in the news
- Read Chapter 1 of **Linguistic Fundamentals for NLP**, if you haven't already.
- Next topic: **supervised learning**. Read:
  - **LxMLS notes** on linear algebra, probability, and classification
  - Optional: survey on word sense disambiguation