

CS5740: Natural Language Processing

Introduction

Instructor: Yoav Artzi

TAs: Rishi Bommasani and Max Grusky

Technicalities

- People:
 - Instructor: Yoav Artzi
 - TAs: Rishi Bommasani and Max Grusky
- Webpage (everything is there):
 - <http://www.cs.cornell.edu/courses/cs5740/2020sp/>
- Discussion group linked from website
- Assignments on CMS
 - Repositories on Github Classroom

Technicalities

- Office hours posted on website
- Example assignments and report layout posted on CMS
- Enrollment
 - Not formally enrolled → email me your NetID for quizzes and CMS ASAP
 - If you are in the system, you should see the course in CMS, if not, you are not in the system
 - This will continue until enrollment stabilizes (about two weeks)
- Slides always posted after class, often appended to the most recent deck lecture

Procedurals

<https://www.cs.cornell.edu/courses/cs5740/2020sp/procedurals.html>

Quizzes

- It is not possible to re-take a missed quiz
- A missed quiz gets zero
- Just like an exam: no copying, chatting, and not taking the quiz remotely → all AI violations
- Come on time
 - Late? Enter quietly and sit at the back
 - Quiz starts on time
- Quiz practice
 - Phones, tablets, or laptops
 - <http://socrative.com>
 - Select “Student Login”
 - Today’s room: NLP2045
 - Use NetID to identify
- After the quiz: please put your electronics aside

Tips

- Work together with your partner, don't simply divide the work
- Discuss with each other
 - Beyond your group
 - This is what the forum is for!

What is this class?

- Depth-first technical NLP course
- Learn the language of natural language processing
- What this class is not?
 - It is not a tutorial to NLTK, PyTorch, etc.
 - There are many resources online that already do that well

Class Goals

- Learn about the issues and techniques of modern NLP
- Be able to read current research papers
- Build realistic NLP tools
- Understand the limitation of current techniques

Main Themes

- Linguistic Issues
 - What is the range of language phenomena?
 - What are the knowledge sources that let us make decisions?
 - What representations are appropriate?
- Modeling and Learning Methods
 - Increasingly complex model structures
 - Learning and parameter estimation
 - Efficient inference
- Engineering Methods
 - Issues of scale
 - Where the theory breaks down (and what to do about it)
- We'll focus on what makes the problems hard, and what works in practice

Three Types of Models

- Generative Models
- Discriminative Models
 - Neural Networks
- Graphical Models

What is NLP?



- Fundamental goal: deep understanding of broad language
 - Not just string processing or keyword matching!
- End systems that we want to build:
 - Simple: spelling correction, text categorization...
 - Complex: speech recognition, machine translation, information extraction, dialog interfaces, question answering...
 - Unknown: human-level comprehension (is this just NLP?)

Today

- State of the art examples
- Some history
- Key problems

Machine Translation

lrytas.lt

The image displays three overlapping screenshots of a news article from the website lrytas.lt, illustrating the evolution of machine translation over time. The background of the article is a lottery ticket with numbers and a large 'X' mark.

2018: The top screenshot shows the article title "The 70-year-old woman 'by mistake' won almost 2 million euros" and the source "dpa-ELTA Inf." dated "2018-01-24 13:33". The text is partially obscured by the other screenshots.

2019: The middle screenshot shows the same article title and source, but with a different layout and a "Comments" button visible.

2020: The bottom screenshot shows the article title "A 70-year-old woman by mistake won nearly 2 million in a lottery. euros" and the source "dpa-ELTA info" dated "2018-01-24 13:33, Updated 2018-01-25 02:57". The text is more complete and includes a "Share" button.

The article text in the 2020 version reads:

In Germany, a [pensioner](#) accidentally noted on a lottery ticket that he was taking part in a Wednesday game and not, as usual, on Saturday, the state-run lottery company said.

The "mistake" turns out to be a success: for six guessing figures, a woman has suffered 1.9 million. euros.

In retrospect, her [retirement](#) commented on her achievement : "Incredible!" She and her husband now want to fulfill their dream of taking a Caribbean cruise.

Personal Assistants



Dialogue Systems

Dialogue Systems

Barbie: Hey, new question. Do you have any sisters?

Tiara: Yeah, I only have one.

Barbie: What's something nice that your sister does for you?

Tiara (tensely): She does nothing nice to me.

Barbie: Well, what is the last nice thing your sister did?

Tiara: She helped me with my project — and then she *destroyed* it.

Barbie: Oh, yeah, tell me more!

Tiara: That's it, Barbie.

Barbie: Have you told your sister lately how cool she is?

Tiara (gritting her teeth): No. She is *not* cool.

Barbie: You never know, she might appreciate hearing it.



Break It!



- Pick an NLP app you can get to now
 - Such as an online translation service or the voice assistant on your phone
- Find examples where it fails
- Try to understand / hypothesize why
 - Hint: it helps to try similar inputs to find the fault lines

Exciting Times

- We can do a lot of things
- But:
 - Must think beyond the NLP technique or the cool system we build
 - Don't forget what we still don't know

NLP History: Pre-statistics

- (1) Colorless green ideas sleep furiously.
- (2) Furiously sleep ideas green colorless

NLP History: Pre-statistics

(1) Colorless green ideas sleep furiously.

(2) Furiously sleep ideas green colorless

It is fair to assume that neither sentence (1) nor (2) (nor indeed any part of these sentences) had ever occurred in an English discourse. Hence, in any statistical model for grammaticality, these sentences will be ruled out on identical grounds as equally "remote" from English. Yet (1), though nonsensical, is grammatical, while (2) is not."
(Chomsky 1957)

NLP History: Pre-statistics

- 70s and 80s: more linguistic focus
 - Emphasis on deeper models, syntax and semantics
 - Toy domains / manually engineered systems
 - Weak empirical evaluation

NLP History: ML and Empiricism

“Whenever I fire a linguist our system performance improves.” –Jelinek, 1988

- 1990s: Empirical Revolution
 - Corpus-based methods produce the first widely used tools
 - Deep linguistic analysis often traded for robust approximations
 - *Empirical evaluation* is essential

NLP History: ML and Empiricism

“Whenever I fire a linguist our system performance improves.” –Jelinek, 1988

- 1990s: Empirical Revolution
 - Corpus-based methods produce the first widely used tools
 - Deep linguistic analysis often traded for robust approximations
 - *Empirical evaluation* is essential

“Of course, we must not go overboard and mistakenly conclude that the successes of statistical NLP render linguistics irrelevant (rash statements to this effect have been made in the past, e.g., the notorious remark, “Every time I fire a linguist, my performance goes up”). The information and insight that linguists, psychologists, and others have gathered about language is invaluable in creating high-performance broad-domain language understanding systems ...”

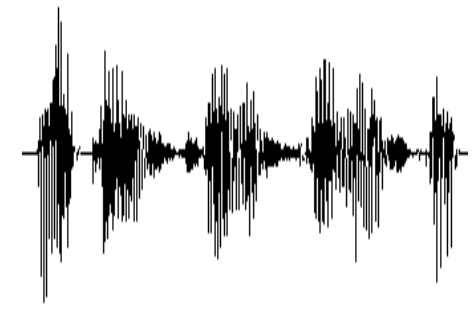
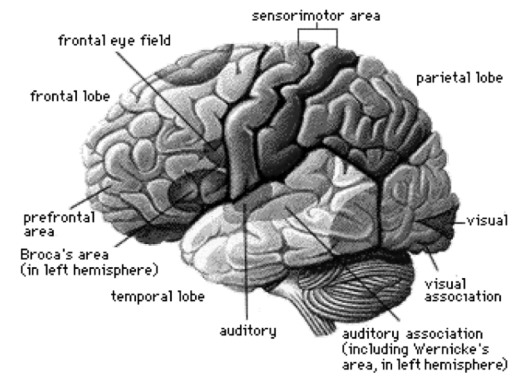
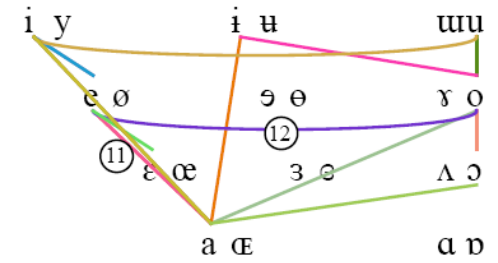
NLP History: ML and Empiricism

“Whenever I fire a linguist our system performance improves.” –Jelinek, 1988

- 1990s: Empirical Revolution
 - Corpus-based methods produce the first widely used tools
 - Deep linguistic analysis often traded for robust approximations
 - *Empirical evaluation* is essential
- 2000s: Richer linguistic representations used in statistical approaches, scale to more data!
- 2010s: NLP+X, excitement about neural networks (again), pre-trained representations
- 2020s: ...

Related Fields

- Computational Linguistics
 - Using computational methods to learn more about how language works
 - We end up doing this and using it
- Cognitive Science
 - Figuring out how the human brain works
 - Includes the bits that do language
 - Humans: the only working NLP prototype!
- Speech
 - Mapping audio signals to text
 - Traditionally separate from NLP, converging?
 - Two components: acoustic models and language models
 - Language models in the domain of stat NLP



Key Problems

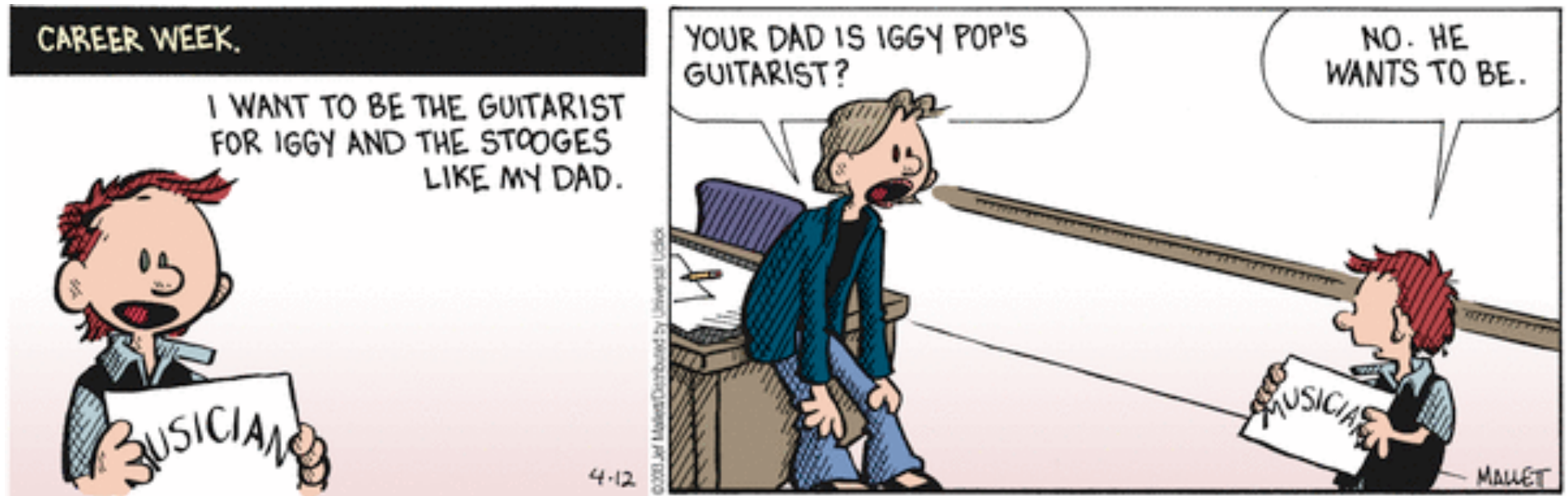
We can understand programming languages.
Why is NLP not solved?

Key Problems

We can understand programming languages.
Why is NLP not solved?

- Ambiguity
- Scale
- Sparsity

Key Problem: Ambiguity

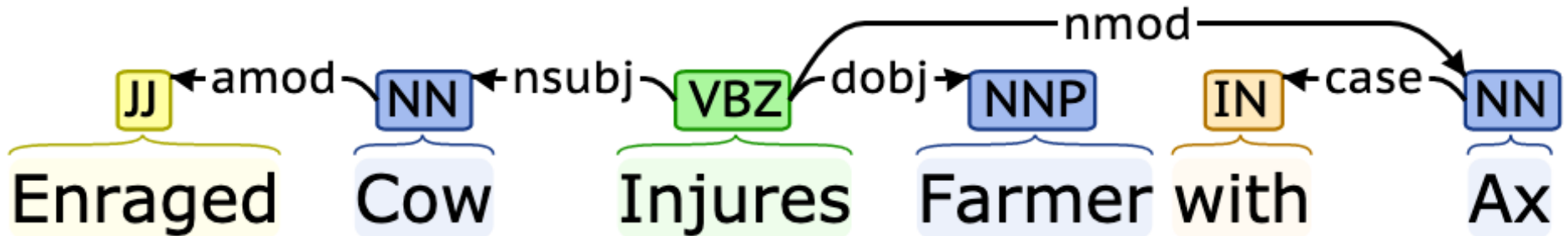


Key Problem: Ambiguity

- Some headlines:
 - Enraged Cow Injures Farmer with Ax
 - Ban on Nude Dancing on Governor's Desk
 - Teacher Strikes Idle Kids
 - Hospitals Are Sued by 7 Foot Doctors
 - Iraqi Head Seeks Arms
 - Stolen Painting Found by Tree
 - Kids Make Nutritious Snacks
 - Local HS Dropouts Cut in Half

Syntactic Ambiguity

Enraged Cow Injures Farmer with Ax



- SOTA: ~95% accurate for many languages when given many training examples, some progress in analyzing languages given few or no examples

Semantic Ambiguity

At last, a computer that understands you like your mother.

Semantic Ambiguity

At last, a computer that understands you like your mother.

- Direct meanings:
 - It understands you like your mother (does) [presumably well]
 - It understands (that) you like your mother
- “*mother*” could mean:
 - a woman who has given birth to a child
 - a stringy slimy substance consisting of yeast cells and bacteria; is added to cider or wine to produce vinegar
- Context matters, e.g. what if previous sentence was:
 - Wow, Amazon predicted that you would need to order a big batch of new vinegar brewing ingredients. ☒

Ambiguities in the Wild



Ambiguities in the Wild: Context

The Atlantic

SUBSCRIBE SEARCH MENU

Susan Collins Unveils a Gun-Control Compromise

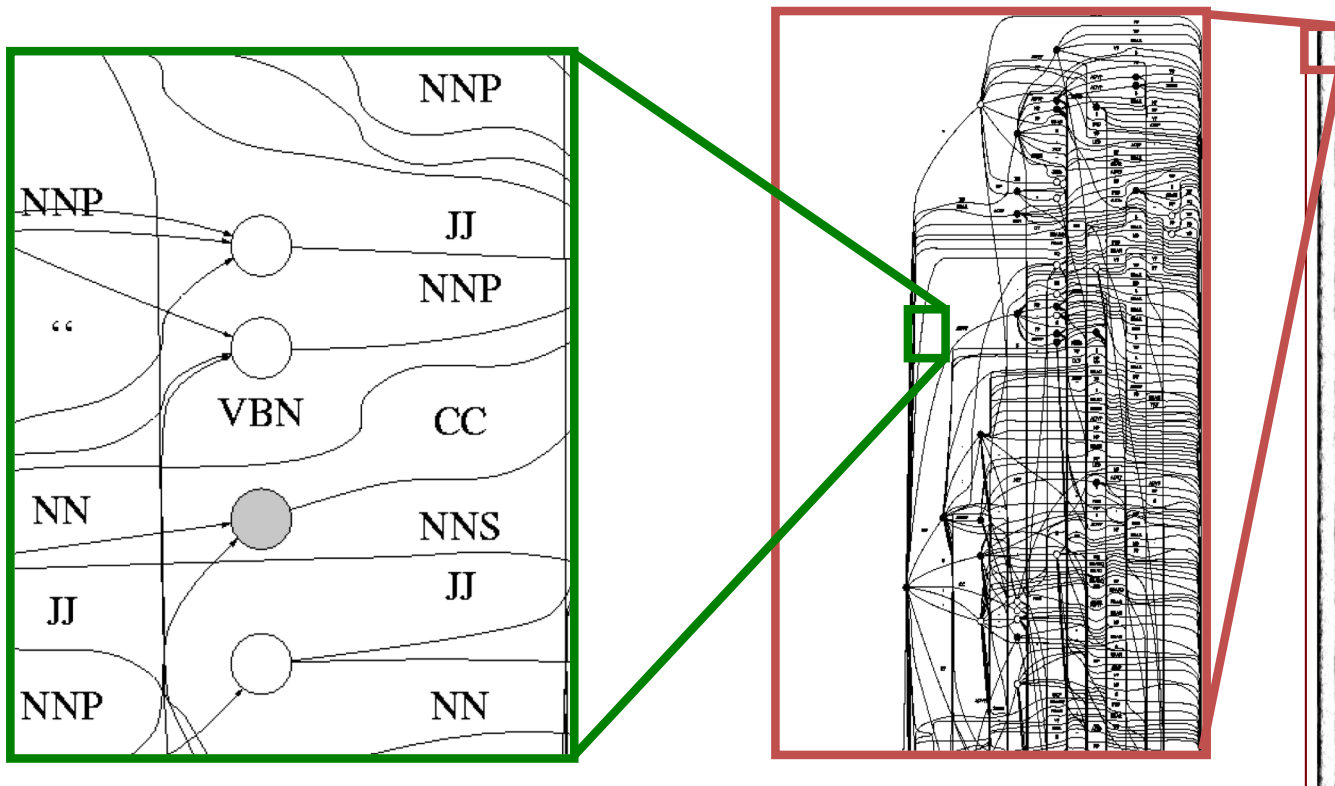
It would restrict sales to individuals on two terrorist watch lists.



Yuri Gripas / Reuters

Key Problem: Scale

- People *did* know that language was ambiguous!
 - ...but they hoped that all interpretations would be “good” ones (or ruled out pragmatically)
 - ...they didn’t realize how bad it would be



Key Problem: Sparsity



- A corpus is a collection of text
 - Often annotated in some way
 - Sometimes just lots of text
 - Balanced vs. uniform corpora
- Examples
 - Newswire collections: 500M+ words
 - Brown corpus: 1M words of tagged “balanced” text
 - Penn Treebank: 1M words of parsed WSJ
 - Canadian Hansards: 10M+ words of aligned French / English sentences
 - The Web: billions of words of who knows what

Key Problem: Sparsity

- However: sparsity is always a problem
 - New unigram (word), bigram (word pair)

