Foundations of Natural Language Processing Lecture 1: Introduction

Mirella Lapata School of Informatics University of Edinburgh mlap@inf.ed.ac.uk



Slides based on content from: Philipp Koehn, Alex Lascarides, Sharon Goldwater, Shay Cohen, Khalil Sima'an, Ivan Titov



- Welcome to Foundations of Natural Language Processing!
- Make sure you are in the right class/room.
- We'll cover course logistics.
- We'll get started on what is NLP and why it is hard.

We assume you are familiar with most/all of the following:

- Basic Python programming
- Finite-state machines, regular languages, context-free grammars
- Dynamic programming (e.g., edit distance, Viterbi, and/or CKY algorithms)
- Concepts from machine learning (e.g., estimating probabilities, making predictions based on data)
- Probability theory (conditional probabilities, Bayes' Rule, independence and conditional independence, expectations)
- Vectors, logarithms, linear algebra, matrix operations
- Some basic linguistic concepts (e.g., parts of speech)

INF2-iads discussed ideas and algorithms for NLP from a largely formal, algorithmic perspective. Here we build on that by:

- Focusing on real data with all its complexities.
- Discussing some of the NLP techniques in more depth.
- Introducing many tasks and technologies that didn't fit into the Inf2-iads story.
- By the end of the course, you'll know how to make your own ChatGPT.

- Course organizer: Ivan Titov
- Lecturers: Mirella Lapata and Ivan Titov
- 3 lectures per week (Tue, Wed & Fri, 12:10–13.00)
- We will use Learn and Drupal for slides, lectures, labs, assignments, due dates, etc
- Tutorials and labs, see here.
- Course discussion forum: Piazza.

Check course website for all the links and up-to-date information

In addition to attending lectures, you are expected to keep up with:

- Readings from textbook: Speech and Language Processing, Jurafsky and Martin: 3rd edition (online) and 2nd edition (paperback, International version, for chapters that aren't updated in 3rd ed).
- There will also be links to academic papers (recommended).
- Tutorials and quizzes.
- Lectures are being recorded. The audience is not in shot.
- Two assignments, worth 25%.
- Exam in April/May, worth 75% of final mark.

What is Natural Language Processing?

	Untitled document @ All suggestions	HEE ASSISTANT 33
English Spanish French French - detected + Paper English Spanish Arabic + Translate	And a space Thee Thee	33 Overall score See performance
Je ne sais past × I do not know!	a cjck email abot a urgnet meeting we ned to hv tomorow/[The meeting will be at 9 am The meeting will be at 9	Goals Adjust goals
	In the <u>bord</u> room. <u>Pis mak sur</u> to <u>krng</u> all the <u>necessary documnts</u> and praper for a	Generative Al 💡
4) ☆ 颵 /	discusion abot the upcoming prject.	All suggostions
iPhone IP	Your You have a set of the	
Feedback List of Presidents of India - Wikipedia, the free encyclopedia enwikipedia.org/wiki.tist.of/ensidents.ofndia - The Prevident of India is the head of state and first citizen of India. The President is also the Commande-in-Cheld of the Indian Armed Forces. Although the Zakir Hussain - Rejendra Prasad - VV Girl - R. Verkataraman Rajendra Prasad - Wikipedia, the free encyclopedia enwikipedia.org/wiki?Bajendra_Prasad - Isten theje info; 3 December 1864 - 28 February 1963) was the first President of the Republic of India A Indian policital duale. Invert y Braina, Prasad	ChatGPT Bin Sity, voli un limerick en français sur Brutus, votre labrador : Un labrador nommé Brutus, tout mignon, Courait après as ballie dans le jardin, sans façon, Avec sa queue en remuement, Il apportait le bonheur, assurément, Brutus, le chien joyeux de la maison. D	

What is Natural Language Processing?

Applications

- Machine Translation
- Information Retrieval
- Question Answering
- Dialogue Systems
- Information Extraction
- Summarization
- Sentiment Analysis
- ...

Core Technologies

- Language modeling
- Part-of-speech tagging
- Syntactic parsing
- Coreference resolution
- Named-entity recognition
- Word sense disambiguation
- Semantic Role Labeling

...

NLP is a big field! We focus mainly on core ideas and methods needed for technologies in the second column (and eventually for applications).

- Core concepts in NLP
- Core linguistic problems and methodologies in NLP
- including machine learning, problem design, and evaluation methods



This is a simple sentence words

Language consists of many levels of structure

- Humans fluently integrate all of these in producing/understanding language
- Ideally, so would a computer!



- Language consists of many levels of structure
- Humans fluently integrate all of these in producing/understanding language
- Ideally, so would a computer!



- Language consists of many levels of structure
- Humans fluently integrate all of these in producing/understanding language
- Ideally, so would a computer!

What does an NLP system need to "know"?



- Language consists of many levels of structure
- Humans fluently integrate all of these in producing/understanding language
- Ideally, so would a computer!

What does an NLP system need to "know"?



- Language consists of many levels of structure
- Humans fluently integrate all of these in producing/understanding language
- Ideally, so would a computer!

What does an NLP system need to "know"?



- Language consists of many levels of structure
- Humans fluently integrate all of these in producing/understanding language
- Ideally, so would a computer!

Levels of Linguistics Analysis



Levels of Linguistics Analysis

Do we really need to model all these levels?



- 1. Who bought a bridge?
- 2. Where will the bridge be re-built?
- 3. How long will it take?

- 1. Who bought a bridge?
- 2. Where will the bridge be re-built?
- 3. How long will it take?

- 1. Who bought a bridge?
- 2. Where will the bridge be re-built?
- 3. How long will it take?

- 1. Who bought a bridge?
- 2. Where will the bridge be re-built?
- 3. How long will it take?

- 1. Who bought a bridge?
- 2. Where will the bridge be re-built?
- 3. How long will it take?

- 1. Who bought a bridge?
- 2. Where will the bridge be re-built?
- 3. How long will it take?

- 1. Who bought a bridge?
- 2. Where will the bridge be re-built?
- 3. How long will it take?

- 1. Who bought a bridge?
- 2. Where will the bridge be re-built?
- 3. How long will it take?

- 1. Who bought a bridge?
- 2. Where will the bridge be re-built?
- 3. How long will it take?

History of NLP





Variability:

He drew the house He made a sketch of the house He showed me his drawing of the house He portrayed the house in his paintings He drafted the house in his sketchbook



Ambiguity:

She drew a picture of herself cart drawn by two horses...

He drew crowds wherever he went ...

The driver slowed as he drew even with me The officer drew a gun and pointed it at ...

- ~ sketched, made a drawing of
- ~ pulled
- ~ attracted
- ~ proceeded
- ~ took out, produced

Why is NLP hard? Ambiguity at many levels

- Homophones: blew and blue
- Word senses: bank (finance or river?)
- Part of speech: chair (noun or verb?)
- Syntactic structure: I saw a girl with a telescope – We'll look into this in more detail!
- Quantifier scope: Every child loves some movie
- Multiple: I saw her duck
- Reference: John dropped the goblet onto the glass table and it broke.
- Discourse: The meeting is canceled. Emily isn't coming to the office today.

How can we model ambiguity, and choose the correct analysis in context?

Syntactic Ambiguity: Prepositional Phrase Attachment



Example with 3 preposition phrases, 5 interpretations:

- Put the block ((in the box on the table) in the kitchen)
- Put the block (in the box (on the table in the kitchen))
- Put ((the block in the box) on the table) in the kitchen
- Put (the block (in the box on the table)) in the kitchen
- Put (the block in the box) (on the table in the kitchen)

The number of parses is an integer series, known as the Catalan numbers!

$$Cat_n = \frac{1}{n+1} \binom{2n}{n} = \frac{(2n)!}{(n+1)!n!}$$

$$C_1 = 1, C_2 = 2, C_3 = 5, C_4 = 14, C_5 = 42, \ldots$$

Syntactic Ambiguity

A typical tree from a standard dataset (Penn treebank WSJ)



Canadian Utilities had 1988 revenue of \$ 1.16 billion , mainly from its natural gas and electric utility businesses in Alberta , where the company serves about 800,000 customers .

Students Cook & Serve Grandparents

On Thursday, September 9, Gorman School hosted the first annual Grandparent's Day.

All Grandparents were invited to a school wide pancake breakfast. Upper grade students served as excellent chefs, as well as taking responsibility for serving the food and the clean up after-

- Ban on Nude Dancing on Governor's Desk
- Iraqi Head Seeks Arms
- Local High School Dropouts Cut in Half
- Red Tape Holds Up New Bridges
- Juvenile Court to Try Shooting Defendant
- Kids Make Nutritious Snacks

Collected by Chris Manning

- Course logistics
- What is Natural Language Processing (AI rock star!)
- Language consists of many levels of structure
- NLP is hard due to ambiguity at many levels

Next lecture: we discuss NLP challenges some more, and probabilistic modeling.