Foundations of Natural Language Processing Lecture 2: Ambiguity and Probabilistic Models

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

- Recap, this is FNLP!
- Continue discussing why NLP is hard
- Why should we use probabilistic models (and machine learning) for NLP?



London Bridge really is falling down. The bridge is being taken apart and moved. Its new home will be a small town in Arizona. This bridge is hundreds of years old. It stretches across the Thames River. Robert McCulloch saw this bridge and decided to bring it to the USA. He paid more than two million dollars. It will cost him more than three million dollars to move it. Each stone will be marked. The pieces must fit when they reach their new home. All that work will not take place overnight. The job will take Six years. The bridge is not small. It is longer than three football fields. It is almost as wide as one football field. In time, the London Bridge will stand high above a new river. Flags will be placed at both ends. Cars will cross it. A small town will be built next to the bridge. Most people in Arizona will never see London. But they will see a part of it in their own state.

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

Why is NLP hard? Ambiguity at many levels



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 ...

- Ambiguity takes place at different levels (syntax, semantics, co-reference, discourse).
- Leads to combinatorial explosion of possible analyses!

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)

Why is NLP hard? Sparsity

any word				
	Туре		Frequency	Tvpe
	the		104.005	NAr
	of		104,325	
	to		92,195	Commissio
	and		66,781	President
	in		62,867	Parliament
	111		57,804	Union
	that		53,683	report
1	S		53,547	Council
	а		45 842	States
	I		-0,0-7L	Oluios

- Let's look at at the frequencies of different words in a large text corpus.
- Assume a "word" is a string of letters separated by spaces
- Most frequent word types in the English Europarl corpus (out of 24M words).

But also, out of 93,638 distinct word types 36,231 (almost 40%) occur only once.

- cornflakes, mathematicians, fuzziness, jumbling
- seudo-rapporteur, lobby-ridden, perfunctorily,
- Lycketoft, UNCITRAL, H-0695
- policyfor, Commissioneris, 145.95, 27a

Order words by frequency. What is the frequency of *n*th ranked word?



To really see what's going on, use logarithmic axes:



Rescaling the Axes



11/27

Zipf's Law

Summarizes the behavior we just saw:

$$f \times r \approx k$$
 or equivalently $f(r) = \frac{k}{r}$

 $\blacksquare f = \text{frequency of a word}$

- **r** = rank of a word (if sorted by frequency)
- k = a constant

Why a line in log-scales?

$$f \times r = k \Rightarrow f = \frac{k}{r} \Rightarrow \log f = \log\left(\frac{k}{r}\right) \Rightarrow \log f = \log k - \log r$$



George Kingsley Zipf (1902–1950)

Regardless of how large our corpus is, there will be a lot of infrequent (and zero-frequency!) words.

- Regardless of how large our corpus is, there will be a lot of infrequent (and zero-frequency!) words.
- In fact, the same holds for many other levels of linguistic structure (e.g., syntactic rules in a CFG).

- Regardless of how large our corpus is, there will be a lot of infrequent (and zero-frequency!) words.
- In fact, the same holds for many other levels of linguistic structure (e.g., syntactic rules in a CFG).
- Why is this a problem?

- Regardless of how large our corpus is, there will be a lot of infrequent (and zero-frequency!) words.
- In fact, the same holds for many other levels of linguistic structure (e.g., syntactic rules in a CFG).

Why is this a problem?

This means we need to find clever ways to estimate probabilities for things we have rarely or never seen during training.

Why is NLP hard? Robustness

Will our NLP systems work across styles, domains, registers, genres?

Suppose we train a part of speech tagger on the Wall Street Journal:

Mr./NNP Vinken/NNP is/VBZ chairman/NN of/IN Elsevier/NNP N.V./NNP ,/, the/DT Dutch/NNP publishing/VBG group/NN ./.

What will happen if we try to use this tagger for social media??

Pikr smh he asked fir yo last name

New words will always be created; NLP algorithms must be robust to new observations.

Why is NLP hard? World Knowledge and Context











Stanford Colors in Context corpus (Monroe et al. 2017)



Stanford Colors in Context corpus (Monroe et al. 2017)



















Natural Language is Grounded



- In this class, we will focus on English
- ... as most of the NLP community does
- However, we would like to develop technologies which work for other languages
- One example: morphology and syntax

Syntactic Diversity



Number of Cases



World Atlas of Language Structures (wals.info)

Number of Cases

Word order freedom and morphology are inter-related. The more freedom in word order:

- The less information is conveyed by word positions
- The more information should be included in the "tokens"
- The richer morphology



Constrained word order Limited or no morphological marking (Relatively) free word order Rich morphology

- variability
- 2 ambiguity
- sparsity
- robustness
- **5** context dependence
- 6 language diversity

Most challenges we discussed can be regarded as manifestations of uncertainty!

Ambiguity: uncertainty with respect to interpretation

Variability: uncertainty in a specific realization for a semantic concept

Robustness: uncertainty with respect to potential inputs

Lack of knowledge (cf sparsity / context dependence): uncertainty! We need: probabilistic models / machine learning Many natural language processing problems can be written mathematically as:

 $\hat{\mathbf{y}} = \operatorname*{argmax}_{\mathbf{y} \in \mathcal{Y}(\mathbf{x})} P(\mathbf{x}, \mathbf{y})$

- **x** is the input, which is an element of a set \mathcal{X}
- y is the output, which is an element of a set \mathcal{Y}
- $P(\mathbf{x}, \mathbf{y})$ is the model, which maps from the set $\mathcal{X} \times \mathcal{Y} \rightarrow [1, 0]$;
- $\hat{\mathbf{y}}$ is the predicted output, which is chosen to maximize $P(\mathbf{x}, \mathbf{y})$.
- finding $\hat{\mathbf{y}}$ is often a search problem
- **e**stimating $P(\mathbf{x}, \mathbf{y})$ is a **learning** problem.

Sketch of a probabilistic model: Spam Classification

- Input \mathcal{X} is a set of emails.
- Output \mathcal{Y} is 1 (spam) or 0 (not spam).
- Model sees lots of (x, y) pairs (emails labeled with 0/1)
- Goal is to learn to predict labels of new, future emails.
- How do we obtain probabilities $P(\mathbf{x}, \mathbf{y})$?
- How do we measure success?

Dear Sir.

First, I must solicit your confidence in this transaction, this is by virture of its nature as being utterly confidencial and top secret. ...

TO BE REMOVED FROM FUTURE MAILINGS, SIMPLY REPLY TO THIS MESSAGE AND PUT "REMOVE" IN THE SUBJECT.

99 MILLION EMAIL ADDRESSES FOR ONLY \$99

Ok, Iknow this is blatantly OT but I'm beginning to go insane. Had an old Dell Dimension XPS sitting in the corner and decided to put it to use, I know it was working pre being stuck in the corner, but when I plugged it in, hit the power nothing happened.

Sketch of a probabilistic model: Spam Classification

- Input \mathcal{X} is a set of emails.
- Output \mathcal{Y} is 1 (spam) or 0 (not spam).
- Model sees lots of (x, y) pairs (emails labeled with 0/1)
- Goal is to learn to predict labels of new, future emails.
- How do we obtain probabilities $P(\mathbf{x}, \mathbf{y})$?
- How do we measure success?

Dear Sir.

first, I must solicit your confidence in this transaction, this is by virture of its nature as being utterly confidencial and top secret. ...



TO BE REMOVED FROM FUTURE AILINGS, SIMPLY REPLY TO THIS MESSAGE AND PUT "REMOVE" IN THE SUBJECT.

99 MILLION EMAIL ADDRESSES FOR ONLY \$99



Ok, Iknow this is blatantly OT but I'm beginning to go insane. Had an old Dell pimension XPS sitting in the corner and decided to put it to use, I know it was working pre being stuck in the corner, but when I plugged it in, hit the power nothing happened.

- There are many challenges which make NLP hard
- Ambiguity is (?) the most fundamental one
 - happens at many levels
 - can lead to a combinatorial explosion in a number 'interpretations'
- Probabilistic modeling is a way to deal with many of these challenges
- Requires decisions at different levels (probability models, algorithms, ...)

Next lecture: we discuss corpora and evaluation methods.