

Natural Language Processing

- Class 01:
 - Introduction
 - Non-randomness of language



16 February 2018
Dr. Katrien Beuls

Natural Language Processing

- Lecturer: Katrien Beuls (katrien.beuls@vub.be)
- Assistants: Jens Nevens
Paul Van Eecke
- Guest lectures by:
 - Luc Steels (Francqui chair)
- Classes:
 - Fridays, 4pm-6pm
 - Room PL9 3.31
- Slides: At the start of each class on PointCarré
- Evaluation
 - Three assignments during semester (50%)
 - Final written exam (50%)
- Check PointCarré regularly for up-to-date info!

Course materials

- Handbook:
Daniel Jurafsky and James H. Martin (2006, 2nd edition). *Speech and language processing*. Prentice Hall.
Third edition (draft):
<https://web.stanford.edu/~jurafsky/slp3/>
- Slides: At the start of each class on PointCarré
- Articles: As PDFs on PointCarré
- Recommended for self-study
 - Christopher Manning and Hinrich Schütze (1999). *Foundations of statistical natural language processing*. MIT Press.
 - Steven Bird, Ewan Klein and Edward Loper (2009). *Natural language processing with Python*. O'reilly.
+free software: the NLP toolkit (www.nltk.org/book/)

What is Natural Language Processing?



- Can we build useful applications that involve natural language input and output?
- Not necessarily by mimicking the brain, but rather through practical **engineering**.

Natural Language Understanding

VS

Understanding Natural Language

- Can we understand how the human mind as a biological mechanism is able to process language?
- Computational modelling as part of **Cognitive Science**

- Can we build useful applications that involve natural language input and output?
- Not necessarily by mimicking the human brain, but rather through practical **engineering**.

Natural Language Understanding

VS

Understanding Natural Language

- Can we understand how the human mind as a biological mechanism is able to process language?
- Computational modelling as part of Cognitive Science

- Can we build useful applications that involve natural language input and output?
- Not necessarily by mimicking the human brain, but rather through practical **engineering**.

Natural Language Understanding

Some examples of applications...

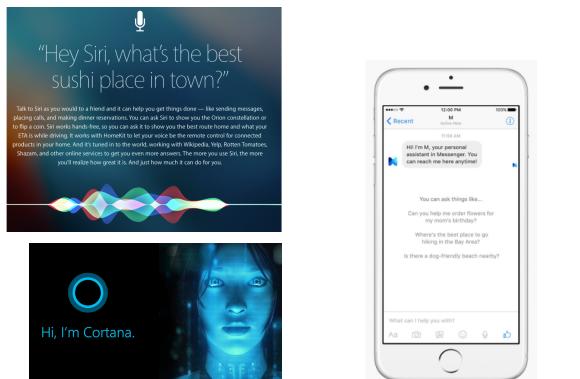
February 2010. Jeopardy!
IBM Watson beats best human player



<https://www.youtube.com/watch?v=P18EdAKuC1U>

Virtual assistants based on NLP

(<https://www.wired.com/2015/08/facebook-launches-m-new-kind-virtual-assistant/>)



Content analytics on big data

ELECTION TRACKER 16
Big Data Analysis of US Election Coverage
Analytics powered by OPINIONTEXT

You are viewing data from the: Last 7 Days

Key Stats		What's Trending?	
No. of articles	8,159	Most popular keywords:	nuclear, abortion, delegate
No. of positive articles	2,477	Visualize >	
No. of negative articles	4,971		
No. of neutral articles	711	Most popular state:	New York
		Visualize >	

Did you know?
Bernie Sanders' March soared in like a lion but is going out like a lamb. He saw 40% of his press trending toward the positive at the beginning of the month but now sees news with his name and a positive tone slipping to 19%.

<http://www.predictiveanalyticstoday.com/election-tracker-big-data-content-analytics-opintext/>

Chatbots relying on NLP

Zuckerberg announces Messenger Platform chatbots

The chief executive of Facebook, Mark Zuckerberg, said outside companies would be able to use **chatbots** on Facebook's Messenger app. (NYTimes April 12, 2016)

Banks Bet on the Next Big Thing: Financial Chatbots

really wants to chat. This week Bank of America, MasterCard and several financial start-ups announced new tools — known as **chatbots** — that will allow customers to ask questions about their financial accounts, initiate (NYTimes October 26, 2016)

Machine translation



Speech-to-speech machine translation

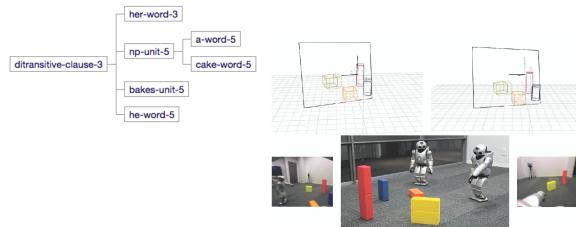


Switch from statistical MT to neural MT

Going Neural (NYTimes December 14, 2016)
Apparently Google Translate, the company's popular machine-translation service, had suddenly and almost immeasurably improved.

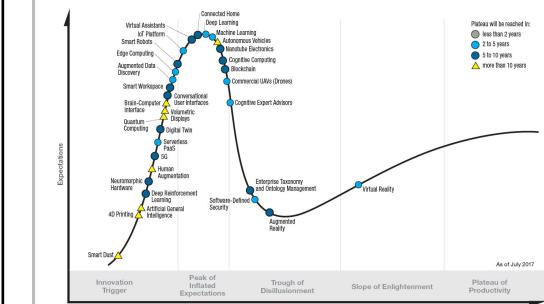
Interactive Language Learning by Robots

Fluid Construction Grammar



<https://www.fcg-net.org/>

Gartner Hype Cycle for Emerging Technologies, 2017



NLP start-ups overviews:

<https://www.crunchbase.com/category/natural-language-processing/789bbbeefc46e1532a68df1f17da87090ea>

<https://angel.co/natural-language-processing>

Why is NLP so difficult?

Three main problems

Ambiguity
Ambiguity
Ambiguity

Rules, but many exceptions
No clear understanding of how humans process language

Examples of ambiguity

Crash Blossoms
Headlines gone wrong



the guardian

home > world > americas > asia

Rio de Janeiro

Mutilated body washes up on Rio beach to be used for Olympics beach volleyball

By Philip Sean Curran Staff Writer

Princess Diana dresses to be auctioned

Ten dresses, including gowns designed by Catherine Walker, Bruce Oldfield and Zandra Rhodes, to go under the hammer

Lauren Cochrane

guardian.co.uk, Friday 15 March 2013 17.26 GMT

Word sense ambiguity

I made her duck.

- I cooked for her (duck=food, her=beneficiary)
- I cooked the duck that belonged to her for dinner (her=owner)
- I caused her to move downwards (duck=action)
- I created a duck for her (duck=artefact)
- I magically turned her into a duck (made=spell)

Ambiguity is pervasive

I caused her to quickly lower her head or body.

- Lexical category: "duck" can be a N or V

I cooked waterfowl belonging to her.

- Lexical category: "her" can be a possessive ("of her") or dative ("for her") pronoun

I made the (plaster) duck statue she owns.

- Lexical semantics: "make" can mean "create" or "cook"

Other difficulties

non-standard English

Great job @justinbieber! Were SOO PROUD of what youve accomplished! U taught us 2 #never say never & you yourself should never give up either

segmentation issues

the New York-New Haven Railroad
the New York-New Haven Railroad

idioms

dark horse
get cold feet
lose face
throw in the towel

neologisms

unfriend
Retweet
bromance

world knowledge

Mary and Sue are sisters.
Mary and Sue are mothers.

tricky entity names

Where is A Bug's Life playing ...
Let It Be was recorded ...
... a mutation on the for gene ...

Three laws of computational linguistics

Hugo Brandt Corstius (1978)

- Whatever you do, semantics will screw things up.
- Every theory, no matter how explicit it is formulated, will turn out to contain errors when you make a program of it.
- The first 80% of accuracy takes little effort, but further diminishing the gap by half takes double the effort of previous work. (~law of diminished returns)



Hugo Brandt Corstius
Dutch author

Hugo Brandt Corstius was a Dutch author. Known for his achievements both in literature and science. In 1970, he was awarded a PhD on the subject of computational linguistics. He was employed at the Mathematisch Centrum in Amsterdam. Wikipedia

Born: August 29, 1935, Eindhoven, Netherlands
Died: February 28, 2014, Amsterdam, Netherlands
Spouse: Henrikele Smit
Parents: Wilhelmina Wytske Molenaar

What is the state-of-the-art?

mostly solved

Spam detection	OK, let's meet by the big ... Dick too small? Buy VIAGRA ...
Text categorization	Philippines shut down Rangers 2-0 Jobless rate hits two-year low
Part-of-speech (POS) tagging	ADJ ADI NOUN VERB ADV Colorless green ideas sleep furiously.
Named entity recognition (NER)	PERSON ORG LOC Obama met with UW leaders in Detroit ...
Information extraction (IE)	You're invited to our Jungs Dinner on Friday May 27 at 8:30pm in Corduba Hall

making good progress

Sentiment analysis	The pho was authentic and yummy. Walter ignored us for 20 minutes.
Coreference resolution	Obama told Mubarak he shouldn't run again.
Word sense disambiguation (WSD)	I need new batteries for my mouse.
Syntactic parsing	I can see Russia from my house!
Machine translation (MT)	Our specialty is panda fried rice. 我们的专长是熊猫饭

still really hard

Semantic search	people protesting globalisation Search ...demonstrators stormed IMF offices...
Question answering (QA)	Q. What currency is used in China? A. The yuan
Textual inference & paraphrase	T. Thirteen soldiers lost their lives ... H. Several troops were killed in the ... YES
Summarization	Sheen continues non against ... Sheen is nuts
Discourse & dialog	Where is Thor playing in SF? Metreon at 4:30 and 7:30

https://www.aclweb.org/aclwiki/index.php?title=State_of_the_art

Many challenges ahead...

• "By ignoring [...] finer details, our language-processing systems have been stuck in an "idiot savant" stage where they can find everything but cannot understand anything. The main language processing challenge of the coming decade is to create robust, accurate, efficient methods that learn to understand the main entities and concepts discussed in any text, and the main claims made."

(Fernando Pereira, Google Research Director, at the META-FORUM 2012)

Watson Doesn't Know It Won on 'Jeopardy'!

IBM invented an ingenious program—not a computer that can think. By JOHN SEARLE (Wall Street Journal Feb. 23, 2011)

What will you learn in this course?



You will learn how to...

- Evaluate which problems need to be solved for a given NLP task;
- Choose an appropriate representation language for solving the task;
- Develop a model in that language;
- Perform inferences with the model;
- Evaluate the model's performance

Foundations

<ul style="list-style-type: none"> • Formal grammars ("Chomsky hierarchy") <ul style="list-style-type: none"> • Type 0 unrestricted grammars • Type 1 context-sensitive grammars • Type 2 context-free grammars • Type 3 regular grammars 	<ul style="list-style-type: none"> • Automata <ul style="list-style-type: none"> • Turing machines • Linearly-bounded automata • Pushdown automata • Finite-state automata (and probabilistic counterparts) 	<ul style="list-style-type: none"> • Shannon Information Theory <ul style="list-style-type: none"> • Perplexity • Entropy • Noisy channel
<ul style="list-style-type: none"> • Linear Algebra <ul style="list-style-type: none"> • Vector representation • Matrix calculus • Singular value decomposition 	<ul style="list-style-type: none"> • Bayesian statistics <ul style="list-style-type: none"> • Inference • Prior distributions • EM-algorithm • Gibbs sampling 	<ul style="list-style-type: none"> • Neural networks <ul style="list-style-type: none"> • Feed-Forward networks • Recurrent Recursive networks • Long Short-Term Memory

Preliminary Schedule

Date	Class	Lecturer	Linguistic phenomenon	techniques	Applications	Assignments
16/02/2018	1Beuls		Introduction / non-randomness	frequency distributions; bigrams & collocations	Overview of applications	
23/02/2018	2Beuls		morphology	regular expressions, finite state automata, string edits	tokenization, lemmatization, spelling correction	
02/03/2018	3Beuls		n-grams	language models, entropy, perplexity	language guessing, Assignment 1	
09/03/2018	4Van Eecke		Parts of speech	Hidden Markov Models	Labelling tasks	
16/03/2018	5Beuls		syntax	Dynamic programming	Dependency parsing	
23/03/2018	6Beuls		Fluid Construction Grammar	Feature structures, unification	Semantic parsing	
30/03/2018	7Beuls		Pragmatics	Inferencing	Language grounding	Assignment 2
06/04/2018	Easter		break			
13/04/2018	Easter		break			
20/04/2018	8Beuls		Document semantics	Vector space models, topic models	Information retrieval	
27/04/2018	9Beuls		Lexical semantics	word spaces, word embeddings	thesaurus extraction, lexical substitution	
04/05/2018	10Beuls		compositional semantics	semantic primitives, distributional semantics	summarization, logical inference	
11/05/2018	11Beuls		opinion mining	sentiment analysis, paraphrase detection	content analytics	Assignment3
18/05/2018	12Beuls		Machine translation	rule-based, statistical, neural MT	general vs domain-specific MT	
25/05/2018	13Beuls		reserve class/ question time			

Levels of linguistic description and corresponding methods



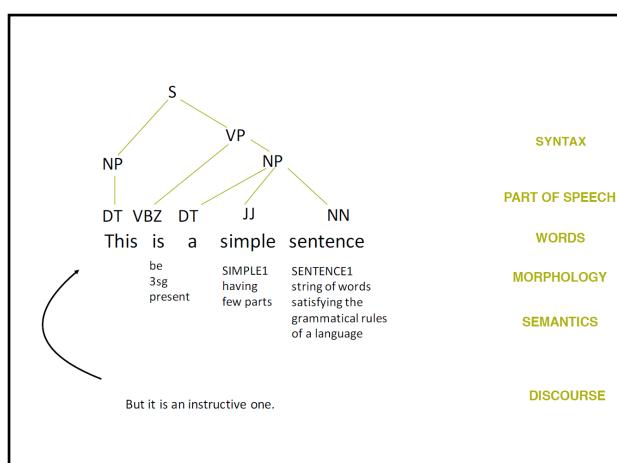
Language Processing

Linguistics

words
morphology
parts of speech
syntax
senses
semantics
discourse

Methods

finite state models
hidden Markov models
grammars and parsing
statistical parsing
sense disambiguation



Non-randomness of Natural Language

Non-randomness 1

What are the relevant units in Natural Language?

Primary data in Natural Language: How can we find recurrent patterns

- **Spoken:** Sound waves, amplitudes varying across time
 - Acoustic signal processing
 - Transcription of individual phonemes
- **Written:** Strings of characters
 - Words: strings separated by spaces, but...
 - Latin: HORVMOMNIVMFORTISSIMISVNTBELGAEPROPTEREAQVODACVL TV ATQUEHUMANITATEPROVINCIALONGISSIMEABSVNT
 - Chinese: 伦敦每日快报指出,两台记载黛安娜王妃一九九七年巴黎
 - Punctuation: >,?!;:
 - Hyphens: sugar-free / sugarfree / sugar free
 - Clitics: Tom's / Kris' / Kris's
 - Compounds: esp. in languages like Dutch & German Rindfleischettikettierungsüberwachungsaufgabenübertragungsgesetz



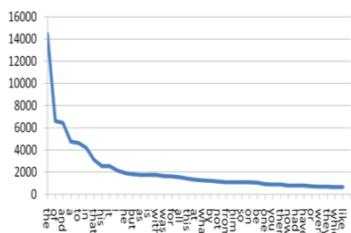
Non-randomness 2

How repetitive is natural language?

Word counts:
Not every word is equally frequent!

	any word	nouns
EUROPARL CORPUS	Frequency in text	Frequency in text
Transcriptions of debates in the European parliament	Token	Token
	1,929,379	129,851
	,	European
	.	Mr
	of	commission
Most frequent words are function words	to	president
They do not convey meaning but functioning as the 'glue' in a language	and	parliament
	in	union
	that	report
	is	council
	a	states
	424,895	member
	424,552	
	Contrast with lexically full words , like nouns	

Only a few words occur very frequently..
...many words occur only once.



- Long right tail of distribution
- 33,447 words occur once, for instance
 - cornflakes
 - mathematicians

Zipf's law

$$f \times r = k$$

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

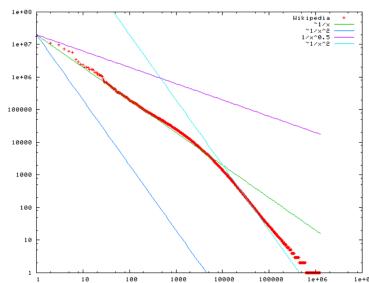
Frequency * frequency-rank:
Linear relation on log-log scale

Zipf's law as
a graph

Linear relationship

$$\begin{aligned} f &\times r = k \\ f &= k/r \\ \log f &= \log k - \log r \end{aligned}$$

=> "power law"



Non-randomness 3

*How do words co-occur
with each other?*

Bigrams: pairs of words

N-grams: sequence of n-words

- Most frequent bigrams are not very interesting: simply combinations of high frequent words (function words)

ALTERNATIVE

- Look at bigrams that are combinations of lexically full words (nouns, adjectives, verbs)
=> multi-word expressions
- Look at bigrams that are combinations of words that co-occur more often than expected by chance => collocations

$C(w^1 w^2)$	w^1	w^2	
80871	of	the	
58841	in	the	
26430	to	the	
21842	on	the	
21839	for	the	
18568	and	the	
16121	that	the	
15630	at	the	
15494	to	be	
13899	in	a	
13689	of	a	
13361	by	the	
13183	with	the	
12622	from	the	
11428	New	York	
10007	he	said	
9775	as	a	
9231	is	a	
8753	has	been	
8573	for	a	

Combinations of lexically full words

Tag Pattern	Example	$C(w^1 w^2)$	w^1	w^2	tag pattern
A N	linear function	7261	United	States	A N
N N	regression coefficients	5412	Los	Angeles	N N
A A N	Gaussian random variable	3301	last	year	A N
A A N	cumulative distribution function	3191	Saudi	Arabia	N N
N A N	mean squared error	2699	last	week	A N
N N N	class probability function	2514	vice	president	A N
N P N	degrees of freedom	2378	Persian	Gulf	A N
		2161	San	Francisco	N N
		2106	President	Bush	N N
		2001	Middle	East	A N
		1942	Saddam	Hussein	N N
		1867	Soviet	Union	A N
		1850	White	House	A N
		1633	United	Nations	A N
		1337	York	City	N N
		1328	oil	prices	N N
		1210	next	year	A N
		1074	chief	executive	A N
		1073	real	estate	A N

Multi-word-expressions:

Combination of words that refer to

- a single entity in reality (Named Entities)
- a single concept that is thought of as a whole, "unit status" ≈ compounds

Collocation: Test if 2 words co-occur more often than expected by chance

Null-hypothesis:

$$\bullet \quad P(w_1 w_2) = P(w_1)P(w_2)$$

Test/quantify deviance

- T-test, Chi-square, log LLH
- Pointwise mutual information

Words can co-occur above chance even if they are not right next to each other

- He knocked yesterday afternoon at my door
- Window-based co-occurrence frequencies

t	$C(w^1)$	$C(w^2)$	$C(w^1 w^2)$	w^1	w^2
4.4721	42	20	20	Ayatollah	Rubollah
4.4721	41	27	20	Bette	Midler
4.4720	30	117	20	Agatha	Christie
4.4720	77	59	20	videocassette	recorder
4.4720	24	320	20	unsalted	butter
2.2446	144	144	20	first	purple
2.2446	13484	10370	20	ever	many
1.3685	14734	13478	20	into	them
1.2176	14093	14776	20	like	people
0.8036	15019	15629	20	time	last

Window 10 left, 10 right

Cooccurrence count	Candidate count	T-score
door	803	27,997 28,305
knock	81	8,191 9,997
doors	98	4,407 9,980
unconscious	47	1,136 6,848
off	545	67,192 23,032
knocked	47	1,136 6,848
down	664	90,098 25,631
sideways	30	924 5,471
knocking	27	915 5,359
opened	75	10,825 8,611
stuffing	23	297 4,793
wall	68	9,971 8,199

MWE and collocations:

More than sum of parts

- **Idiomatic:** conventional way of saying things
 - Let's have a cup of strong tea
- **Non compositionality:** meaning of the expression cannot be fully predicted from its parts
 - Collocations put constraints on meaning: strong tea taste
 - Difficult to substitute with near-synonyms powerful tea
 - MWEs refer to concepts, units: e.g. real estate
 - Idioms are extreme examples of non-compositionality
 - He kicked the bucket
 - She was caught between two stools
 - It cost me an arm and a leg
 - Not freely modifiable He kicked the water bucket??
 - Not literally translatable Il a lancé le seau??

NLP with Python Natural Language ToolKit (NLTK)

Checklist

- Texts are represented in Python using lists: ['Monty', 'Python']. We can use indexing, slicing, and the len() function on lists.
- A word "token" is a particular appearance of a given word in a text; a word "type" is the unique form of the word as a particular sequence of letters. We count word tokens using len(text) and word types using len(set(text)).
- We obtain the vocabulary of a text t using sorted(set(t)).
- To derive the vocabulary, collapsing case distinctions and ignoring punctuation, we can write set(w.lower() for w in text if w.isalpha()).
- A frequency distribution is a collection of items along with their frequency counts (e.g., the words of a text and their frequency of appearance).

What is NLTK?

Language processing task	NLTK modules	Functionality
Accessing corpora	corpus	standardized interfaces to corpora and lexicons
String processing	tokenize, stem	tokenizers, sentence tokenizers, stemmers
Collocation discovery	collocations	t-test, chi-squared, point-wise mutual information
Part-of-speech tagging	tag	n-gram, backoff, Brill, HMM, TrfT
Machine learning	classify, cluster, tbl	decision tree, maxent, entropy, naive Bayes, EM, k-means
Chunking	chunk	regular expression, n-grams, entity
Paraphrases	parse_cog	chart, feature-based, unification, probabilistic, dependency
Semantic interpretation	sem, inference	lambda calculus, first-order logic, model checking
Evaluation metrics	metrics	precision, recall, agreement coefficients
Probability and estimation	probability	frequency distributions, smoothed probability distributions
Applications	app, chat	graphical concordancer, parsers, WordNet browser, chatbots
Linguistic fieldwork	toolbox	manipulate data in SIL Toolbox format

Homework

- Install NLTK
- Go through commands and exercises in Chapter 1
- <http://www.nltk.org/book/ch01.html>

Next week

Date	Class	Lecturer	Linguistic phenomenon	techniques	Applications	Assignments
16/02/2018	1	Beuls	introduction / non-randomness	Frequency Distributions: words & collocations	Overview of applications	
23/02/2018	2	Beuls	morphology	regular expressions, finite state automata, n-grams	tokenization, lemmatization, morphological analysis	
02/03/2018	3	Beuls	n-grams	language models, entropy, perplexity	language guessing,	Assignment 1
09/03/2018	4	van Eecke	Parts of speech	Hidden Markov Models	labelling tasks	
16/03/2018	5	Beuls	syntax	Dynamic programming	Dependency parsing	
23/03/2018	6	Beuls	Fluid Construction Grammar	Feature structures, unification	Semantic parsing	
30/03/2018	7	Beuls	Pragmatics	Inferencing	Language grounding	Assignment 2
06/04/2018	Easter		break			
13/04/2018	Easter		break			
20/04/2018	8	Beuls	Document semantics	Vector space models, topic models	Information retrieval thesaurus extraction, lexical substitution	
27/04/2018	9	Beuls	Lexical semantics	word spaces, word embeddings	formal semantics, distributional semantics	
04/05/2018	10	Beuls	compositional semantics	semantics	summarization, logical inference	
11/05/2018	11	Beuls	opinion mining	sentiment analysis, paraphrase detection	content analytics	Assignment3
18/05/2018	12	Beuls	Machine translation	rule-based, statistical, neural MT	general vs domain-specific MT	
25/05/2018	13	Beuls	reserve class/ question time			