

#### Intro to NLP

Yoav Goldberg

Technion, 2018



#### About me

- Working with NLP / Text / ML for ~10 years.
- Currently: Faculty at Bar Ilan University NLP infrastructure, word similarity, deep learning for language
- Before: Google Research New York NLP infrastructure and toolchain, Machine Translation
- Before: Ben Gurion University (msc + phd) Hebrew NLP, syntax and parsing, structure prediction
- Before that: computer security / cyber / etc.



#### This Course

- Natural Language Processing (NLP):
  - Algorithms to process, analyze and understand texts in natural language.
    - Understanding **Structure**
    - Understanding Meaning
  - Before you solve a problem, you need to understand it.
    - What is natural language? what are the building blocks?



#### This Course

- Natural Language Processing (NLP):
  - Algorithms to process, analyze and understand texts in natural language.
    - Understanding **Structure**
    - Understanding Meaning
  - Applications:
    - Helping people do things with text.
  - Before you solve a problem, you need to *understand* it.
    - What is natural language? what are the building blocks?



#### What

- How to think about working with text data.
- Understand the challenges.
- Understand the linguistic concepts.
- Understand the ML concepts.
- Understand the algorithms.
- Understand where algorithms fail, and how to asses the quality of algorithms.
- Understand what is missing.



#### What Not

- I assume you know basic statistics / ML.
- I will only briefly mention neural nets / deep learning.
- In many cases I will show the lower-level algorithms, without much details on how their outputs is used---unless you ask me to.



#### Formalities

- 3 hours a week, given by me.
- No distinction between "Lecture" and "Tirgul" sessions.
- Four programming assignments (40% of the grade)
- Final exam (rest of the grade).



## Assignments

- Different from each other, covering different levels of NLP work.
- Can be done in Python or Java (I advocate for Python)
- Its good if you are familiar with Linux/Commandline.
- Assignments require also a written report component.

#### Discussion / Communication

- Discussion about assignments will be done through a Piazza forum. Don't email me with questions, everything should go through the piazza. piazza.com/technion.ac.il/spring2018/236299
- Communications with me regarding course, also through piazza (private messages), not via email.



#### Course

- Very similar, but not identical, to the course we teach at Bar-Ilan.
- Last semester's website: <u>http://www.cs.biu.ac.il/~89-680/</u>



#### External Materials

- No books covers the entire class.
- I will post links to related courses/materials on the course website.
- Each class will be accompanied by a list of related papers. You are encouraged to read them.











Google	wha	t is nlp					I C
	All	Videos	Images	News	Maps	More	Settings Tools
	Abou	t 12,400,000	results (0.68	seconds)			

**Neuro-linguistic programming (NLP)** is an approach to communication, personal development, and psychotherapy created by Richard Bandler and John Grinder in California, United States in the 1970s.

Neuro-linguistic programming - Wikipedia https://en.wikipedia.org/wiki/Neuro-linguistic\_programming

About this result . Feedback



oogle	what	t is natural	language	processir	ng		ļ	Q
	All	Videos	Images	News	Maps	More	Settings	Tools
	Abour	t 8,640,000 r	esults (0.69 s	seconds)				

**Natural language processing** is a field of computer science, artificial intelligence, and computational linguistics concerned with the interactions between computers and human (**natural**) languages. As such, **NLP** is related to the area of human– computer interaction.

Natural language processing - Wikipedia https://en.wikipedia.org/wiki/Natural\_language\_processing

**Cill alway** 

About this result . Feedback



Google	10 ar	m in tel avi	v in new ;	york				Ļ	٩
	All	Images	News	Maps	Videos	More	Setting	gs	Tools
	About	7,760,000 re	esults (0.87	seconds)					
	10:00 <b>3:0</b>	0 Thursday 0 Thurs	/, in Tel A sdav. i	viv-Yafo i n New	₅ ⁄York.	NY. USA			

Jarool Time to New York Time Convertor - Time Die









to me 🗘

Did you send it before or after our conversation earlier?







Pierre Vinken, 61 years old, will join the board as a nonexecutive director, Nov. 21, 1987

 $\Downarrow$  chunking









#### sentiment analysis

#### The soup, which I expected to be good, was bad





#### parsing



#### Structure of Sentences



#### parsing helps sentiment



#### October 14, 2002, 4:00 a.m. PT

For years, <u>Microsoft Corporation CEO Bill Gates</u> railed against the economic philosophy of opensource software with Orwellian fervor, denouncing its communal licensing as a "cancer" that stifled technological innovation.

Today, Microsoft claims to "love" the open-source concept, by which software code is made public to encourage improvement and development by outside programmers. Gates himself says Microsoft will gladly disclose its crown jewels--the coveted code behind the Windows operating system--to select customers.

"We can be open source. We love the concept of shared source," said <u>Bill Veghte</u>, a <u>Microsoft VP</u>. "That's a super-important shift for us in terms of code access."

**<u>Richard Stallman</u>**, <u>founder</u> of the <u>Free Software</u> <u>Foundation</u>, countered saying...



#### information extraction

NAME	TITLE	ORGANIZATION
Bill Gates	CEO	Microsoft
Bill Veghte	VP	Microsoft
Richard Stallman	founder	Free Soft



#### NLP In a Nutshell



non-trivial useful output

takes as input text in human language and process it in a way that suggests an intelligent process was involved







"Meaning" "insights"



#### Text in human language



#### Text in another language



structured data in some format



# (some) NLP Applications

- Search
- Knowledge Discovery
- Translation
- Question Answering
- Ads / market analysis
- Chatbots / personal assistants
- Others? (discuss)



# Interdisciplinary Field

- NLP Combines:
  - AI
  - Linguistics
  - Machine Learning / statistics
  - Algorithms (dynamic programming, search)



#### Related Fields

- Computational Linguistics
  - Using algorithms to understand language.
- Psycholinguistics / Cognitive Science
  - How is language processed by the human brain / by humans?
- Social Linguistics
  - How do humans use and shape language as a social construct?



#### How do we do NLP?

- 1950 -- ~1990s ---> Write many rules
- 1990s -- ~2000s ---> Corpus based statistics
- 2000s -- ~2014 ---> Supervised machine learning
- 2014 -- today ---> "deep learning"



#### How do we do NLP?

- 1950 -- ~1990s ---> Write many rules
- 1990s -- ~2000s ---> Corpus based statistics
- 2000s -- ~2014 ---> Supervised machine learning
- 2014 -- today ---> "deep learning"

2021+ --> write rules, aided by machine learning.



#### How do we do NLP?

- 1950 -- ~1990s ---> Write many rules
- 1990s -- ~2000s ---> Corpus based statistics
- 2000s -- ~2014 --->Supervised machine learning
- 2014 -- today ---> "deep learning"

2021+ --> write rules, aided by machine learning.



# "Empirical Revolution"

"Whenever I fire a linguists, performance goes up" (Fred Jelinek, ~1988)

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;** for instance, in the speech recognition setting described above, better understanding of language structure can lead to better language models." (Lilian Lee, 2001)


### Why is language Interesting?

- Invented by humans for communications.
- Learned from experience, in a transparent manner.
- Preserves ideas across space and time.
- Symbolic system.
  - ...also for abstract stuff.
- Language facilitates thinking?
- Tons of human knowledge encoded in language.

- Phonetics (how sounds are made: lips, spit, throat)
- Phonology (how sounds can combine)
- Morphology (how words are built)
- Syntax (how words are combined)
- Semantics (the meaning of words/phrases)
- Pragmatics (the true meaning of words/phrases)
- Discourse (structure/meaning across sentences)

- Phonetics (how sounds are made: lips, spit, throat)
- Phonology (how sounds can combine)
- Morphology (how words are built)
- Syntax (how words are combined)
- Semantics (the meaning of words/phrases)
- Pragmatics (the true meaning of words/phrases)
- Discourse (structure/meaning across sentences)

- Phonetics (how sounds are made: lips, spit, throat)
- Phonology (how sounds can combine)
- Morphology (how words are built)
- Syntax (how words are combined)
- [מצאתי [מטבע [על הרצפה]]]
  Semantics (the meaning of words/phrases)
- Pragmatics (the true meaning of words/phrases)
- Discourse (structure/meaning across sentences)

- Phonetics (how sounds are made: lips, spit, throat)
- Phonology (how sounds can combine)
- Morphology (how words are built)
- Syntax (how words are combined)
- [מצאתי [מטבע [על הרצפה]]] • Semantics (the meaning of words/phrases) ״יש לך שקל?״ -> האם יש ברשותך שקל
- Pragmatics (the true meaning of words/phrases)
- Discourse (structure/meaning across sentences)

- Phonetics (how sounds are made: lips, spit, throat)
- Phonology (how sounds can combine)
- Morphology (how words are built)
- Syntax (how words are combined)
- [מצאתי [מטבע [על הרצפה]]]
  Semantics (the meaning of words/phrases)
  - ״יש לך שקל?״ -> האם יש ברשותך שקל
- Pragmatics (the true meaning of words/phrases)
  "יש לך שקל?" -> תן לי שקל / רוצה שקל?
- Discourse (structure/meaning across sentences)

- Phonetics (how sounds are made: lips, spit, throat)
- Phonology (how sounds can combine)
- Morphology (how words are built)
- Syntax (how words are combined)
- [מצאתי [מטבע [על הרצפה]]] • Semantics (the meaning of words/phrases) ״יש לך שקל?״ -> האם יש ברשותך שקל
- Pragmatics (the true meaning of words/phrases)
  "יש לך שקל?" -> תן לי שקל / רוצה שקל?
- Discourse (structure/meaning across sentences)

#### בגלל, ברם, ואכן...

- Phonetics (how sounds are made: lips, spit, throat)
- Phonology (how sounds can combine)
- Morphology (how words are built)
- Syntax (how words are combined)
- [מצאתי [מטבע [על הרצפה]]]
  Semantics (the meaning of words/phrases)
  - ״יש לך שקל?״ -> האם יש ברשותך שקל
- Pragmatics (the true meaning of words/phrases)
  "יש לך שקל?" -> תן לי שקל / רוצה שקל?
- Discourse (structure/meaning across sentences)

#### בגלל, ברם, ואכן... הוא אמר ש..

# Why is Language Difficult

- Variability
- Ambiguity
- Creativity
- Coverage
- Generalization

# Variability and Ambiguity







### Variability

he acquired it

he purchased it

he bought it

it was bought by him

it was sold to him

she sold it to him

she sold him that



## Variability

2m 2ma 2mar 2mara 2maro 2marrow 2mor 2mora 2moro 2morow 2morr 2morro 2morrow 2moz 2mr 2mro 2mrrw 2mrw 2mw tmmrw tmo tmoro tmorrow tmoz tmr tmro tmrow tmrrow tmrrw tmrw tmrw tmw tomaro tomarow tomarro tomarrow tomm tommarow tommarrow tomarro tommorow tommorrow tommorw tommoro tomoro tomorow tomorro tomorrw tomoz tomrw tomz



bank

apple

star



book

book



#### book a flight

#### read a book



ספר



#### הרכבת



Ambiguity

the book was interesting

the movie was interesting



#### ספר עזר לרופא שיניים בהוצאת כתר



#### שמעתי ש ספר עזר לרופא שיניים בהוצאת כתר קניתי



Ambiguity

I ate pizza with olives

I ate pizza with friends



#### I saw her duck

 Self control is waiting for the movie to start to eat popcorn.



#### נעצר עורך דין שניצל ופיתה שופטת וגנב



Ambiguity



#### וטרינר רצה להחרים טווס - וירה באישה עם מגרפה ynet.co.il



- News headlines
  - Ban on Nude Dancing on Governor's Desk
  - Teacher Strikes Idle Kids
  - Hospitals Are Sued by 7 Foot Doctors
  - Iraqi Head Seeks Arms
  - Kids Make Nutritious Snacks
  - Local Highschool Dropouts Cut in Half



#### חולצה מטיילת בהודו



fat people eat candy



#### fat people eat accumulates

- Phonetics (how sounds are made: lips, spit, throat)
- Phonology (how sounds can combine)
- Morphology (how words are built)
- Syntax (how words are combined)

- There is ambiguity on all levels. (can you think of examples?)
- Semantics (the meaning of words/phrases)
- Pragmatics (the true meaning of words/phrases)
- Discourse (structure/meaning across sentences)





#### Zipf law

- "the frequency of any word is inversely proportional to its rank in the frequency table"
- Word frequencies follow a power-law distribution.
- --> Long tail
- --> Most "events" will occur only few times, if any.





Wei Xu o socialmedia-class.org





 frequency of word is inversely proportional to its rank in the frequency table

word frequency in the Brown corpus







Unique types vs. Million tokens



Million tokens



### Coverage

in a **43M** words text, there are:

- 316,710 unique words
- 144,999 words occur only once
- 42,525 words occur 2 times
- 21,618 words occur 3 times
- 13,306 words occur 4 times
- 9,488 words occur 5 times
- 26,024 words appear >50 times


in a **43M** words text, there are:

- 316,710 unique words
- 144,999 words occur only once
- 42,525 words occur 2 times
- 21,618 words occur 3 times
- 13,306 words occur 4 times
- 9,488 words occur 5 times
- 26,024 words appear >50 times

#### how may words occur 0 times?



Some 1-count words:

- poetically
- unregulation
- PORTLAND
- timbers
- wildlands
- philologist
- Metallurgic
- 8,553
- baby-killer
- sofia-based

Some 5-count words:

- hologram
- 500.7
- Ben-Gurion
- yahoo
- Nineteenth
- profitabilty
- ballpoint
- interstate-highway
- nested
- algorithm



Some 10-count words:

- pivot

- sensitively
- provable
- cathode-ray
- social-democratic
- encyclical
- anti-military
- marinated
- transparently

Some 50-count words:

- gadgets
- treasurers
- narrator
- 349
- Syrians
- Kearns
- drummer
- trophy
- vanishing



- No matter how large the training corpus is:
  - there are very likely to be word forms we did not see
  - there will be many word-pairs we didn't see
  - there will be even more word-triplets we didn't see
  - there will be even more sentences we didn't see



## Word Clustering

mailman	10000011010111
salesman	100000110110000
bookkeeper	1000001101100010
troubleshooter	10000011011000110
bouncer	10000011011000111
technician	1000001101100100
ianitor	1000001101100101
saleswoman	1000001101100110
Nike	101101110010010101011100
Maytag	1011011100100101010111010
Generali	1011011100100101010111011
Gap	10110111001001010111110
Harley-Davidson	1011011100100101010111110
Enfield	10110111001001010101111110
genus	10110111001001010101111111
Microsoft	10110111001001011000
Ventritex	101101110010010110010
Tractebel	1011011100100101100110
Synopsys	1011011100100101100111
WordPerfect	1011011100100101101000
John	101110010000000000
Consuelo	101110010000000000
Jeffrey	10111001000000010
Kenneth	1011100100000001100
Phillip	101110010000000011010
WILLIAM	101110010000000011011
Timothy	10111001000000001110



Creativity

#### ENTERTAINMENT

#### Lindsay Lohan Just Went Full Regina George On Ariana Grande's Selfie

Raise your hand if you've ever felt personally victimized by Lindsay Lohan.

() 11/19/2016 01:27 pm ET | Updated 4 days ago



### Generalization

 Discrete events --> hard to generalize to unseen events.



- Discrete symbols
- Sequences
- Hierarchies



- Discrete symbols
- Sequences
- Hierarchies

- many unseen events
- hard to generalize
- hard to use math
- sparse representations



- Discrete symbols
- Sequences

- order matters
- sub sequences matter

• Hierarchies



- Discrete symbols
- Sequences
- Hierarchies

document --> paragraphs --> sentences
 --> phrases --> words --> morphemes



- Discrete symbols
- Sequences
- Hierarchies

- which level to look at?
- sequences are not enough

document --> paragraphs --> sentences
 --> phrases --> words --> morphemes

# Kinds of NLP Applications

- document --> label
- sentence --> label
- sentence(s) --> sentence(s)
- many documents / many sentences (aggregate)
  --> Information/relations (Information extraction 1)
- document / sentence --> information/relations (Information extraction 2)
- sentence --> API command (semantic parsing)
- clustering



## document --> label

- language classification
- topic classification
- author classification
- sentiment classification
- interestingness
- relevance



### sentence --> label

- mostly the same as in "document --> label" but on a more granular level.
  - (but: shorter text --> harder task)



### sentence(s) --> sentence(s)

- translation
- question answering
- simplification
- summarization
- query completion
- email auto-response generation



- Information extraction (1)
- "Table filling"

"bill gates", "founder of", "microsoft" "bethoven", "born in", "bonn"



#### As a task: Filling slots in a database from sub-segments of text.

October 14, 2002, 4:00 a.m. PT

For years, Microsoft Corporation CEO Bill Gates railed against the economic philosophy of opensource software with Orwellian fervor, denouncing its communal licensing as a "cancer" that stifled technological innovation.

Today, Microsoft claims to "love" the open-source concept, by which software code is made public to encourage improvement and development by outside programmers. Gates himself says Microsoft will gladly disclose its crown jewels--the coveted code behind the Windows operating system--to select customers.

"We can be open source. We love the concept of shared source," said Bill Veghte, a Microsoft VP. "That's a super-important shift for us in terms of code access."

Richard Stallman, founder of the Free Software Foundation, countered saying...





#### As a task: Filling slots in a database from sub-segments of text.

October 14, 2002, 4:00 a.m. PT

For years, <u>Microsoft Corporation CEO Bill Gates</u> railed against the economic philosophy of opensource software with Orwellian fervor, denouncing its communal licensing as a "cancer" that stifled technological innovation.

Today, Microsoft claims to "love" the open-source concept, by which software code is made public to encourage improvement and development by outside programmers. Gates himself says Microsoft will gladly disclose its crown jewels--the coveted code behind the Windows operating system--to select customers.

"We can be open source. We love the concept of shared source," said <u>Bill Veghte</u>, a <u>Microsoft VP</u>. "That's a super-important shift for us in terms of code access."

**<u>Richard Stallman</u>**, <u>founder</u> of the <u>Free Software</u> <u>Foundation</u>, countered saying...



NAME	TTT.E	ORGANTZATION
Bill Gates	CEO	Microsoft
Bill Veghte	VP	Microsoft
Richard Stallman	founder	Free Soft



### As a family of techniques:

Information Extraction =

segmentation + classification + clustering + association

October 14, 2002, 4:00 a.m. PT

For years, <u>Microsoft Corporation CEO Bill Gates</u> railed against the economic philosophy of opensource software with Orwellian fervor, denouncing its communal licensing as a "cancer" that stifled technological innovation.

Today, <u>Microsoft</u> claims to "love" the open-source concept, by which software code is made public to encourage improvement and development by outside programmers. <u>Gates</u> himself says <u>Microsoft</u> will gladly disclose its crown jewels--the coveted code behind the Windows operating system--to select customers.

"We can be open source. We love the concept of shared source," said <u>Bill Veghte</u>, a <u>Microsoft VP</u>. "That's a super-important shift for us in terms of code access."

Richard Stallman, founder of the Free Software Foundation, countered saying...

Microsoft Corporation CEO Bill Gates Microsoft Gates aka "named entity Microsoft extraction" Bill Veghte Microsoft VP Richard Stallman founder Free Software Foundation



### As a family of techniques:

Information Extraction = segmentation + classification + association + clustering

October 14, 2002, 4:00 a.m. PT

For years, <u>Microsoft Corporation CEO Bill Gates</u> railed against the economic philosophy of opensource software with Orwellian fervor, denouncing its communal licensing as a "cancer" that stifled technological innovation.

Today, <u>Microsoft</u> claims to "love" the open-source concept, by which software code is made public to encourage improvement and development by outside programmers. <u>Gates</u> himself says <u>Microsoft</u> will gladly disclose its crown jewels--the coveted code behind the Windows operating system--to select customers.

"We can be open source. We love the concept of shared source," said <u>Bill Veghte</u>, a <u>Microsoft VP</u>. "That's a super-important shift for us in terms of code access."

**<u>Richard Stallman</u>**, <u>founder</u> of the <u>Free Software</u> <u>Foundation</u>, countered saying... Microsoft Corporation CEO Bill Gates Microsoft Gates Microsoft Bill Veghte Microsoft VP Richard Stallman founder Free Software Foundation



### As a family of techniques:

Information Extraction = segmentation + classification + association + clustering

October 14, 2002, 4:00 a.m. PT

For years, Microsoft Corporation CEO Bill Gates	
railed against the economic philosophy of open-	<b>Microsoft Corporation</b>
its communal licensing as a "cancer" that stifled	CEO
technological innovation.	Bill Gates
Today, Microsoft claims to "love" the open-source	Microsoft
concept, by which software code is made public to encourage improvement and development by	Gates
outside programmers. <u>Gates</u> himself says	Microsoft
Microsoft will gladly disclose its crown jewelsthe coveted code behind the Windows operating	Bill Veghte
systemto select customers.	Microsoft
"We can be open source. We love the concept of	VP
shared source," said <u>Bill Veghte</u> , a <u>Microsoft VP</u> .	Richard Stallman
code access."	founder
Pichard Stallman, foundar of the Eree Software	Free Software Foundation
Foundation, countered saying	



### As a family of techniques:

#### Information Extraction = segmentation + clustering

October 14, 2002, 4:00 a.m. PT

For years, <u>Microsoft Corporation CEO Bill Gates</u> railed against the economic philosophy of open- source software with Orwellian fervor, denouncing its communal licensing as a "cancer" that stifled technological innovation.	*	Microsoft Corporation CEO Bill Gates			icrosoft	icrosoft	ree Soft
Today, <u>Microsoft</u> claims to "love" the open-source	*	Microsoft			∑ }	Σ	E4
encourage improvement and development by		Gates					ler
outside programmers. <u>Gates</u> himself says	*	Microsoft		т. т.	0		nno
coveted code behind the Windows operating		Bill Veghte		E	U E O	ΔP	fo
systemto select customers.		Microsoft					G
"We can be open source. We love the concept of	*	VP					lma:
shared source," said <u>Bill Veghte</u> , a <u>Microsoft VP</u> . "That's a super-important shift for us in terms of		Richard Stallman	ī 🖌			Ø	tal
code access."		founder			tes	sght	Ň
Richard Stallman, founder of the Free Software		<b>Free Software Foundation</b>		E	.1 Ga	.1 V€	hard
Foundation, countered saying			J	NAN	Bil	Bil	Ric



# document / sentence --> information/relations

- Information extraction (2)
- similar to IE1, but harder (cannot rely on redundancy)
- "facts about the text" / "what is this about"
- Google's "calendar events from email" feature.

# Sentence --> API command (semantic parsing)

- "busses from tel-aviv to petach-tikva on thursday"
- "show me images of dogs"
- "add a meeting with mark on tuesday afternoon"
- "send sms to wife buy milk"
- "when was bethoven born"



# Clustering

- Find similar texts
- Summarize collections of similar texts



#### the basic units of text processing

Document Collection / Cluster

Document

Section

Paragraph

Sentence

Document Collection / Cluster

Document

Section

by date by topic by web domain

. . .

Paragraph

Sentence

Document Collection / Cluster

Document

Section Subsection Paragraph

Sentence

**Document Collection / Cluster** Document Section Subsection Paragraph Sentence Phrases... Word

**Document Collection / Cluster** Document Section Subsection Paragraph Sentence Phrases... Word

Characters, morphemes

Document Collection / Cluster

#### Document

Section

Paragraph

Sentence



### What is a document?



📄 edition.cnn.com/2016/06/21/opinions/brexit-campaign-disaffection-wismayer... ☆ 保 🜉 🗢

#### Regions + CNN

#### International Edition + menu

Yoav

 $\equiv$ 

British lawmakers clash ahead of Brexit vote 05:57

#### Story highlights

Britons will vote Thursday on whether their country should remain in the European Union Editor's Note: Henry Wismayer is a writer and commentator based in London. You can follow him @henrywismayer. The views expressed are his own.

Henry Wismayer: Urge to leave has caught the wind of native disaffection

(CNN) — For any external observer curious to understand the issues shaping the debate over Britain's membership of

#### Top stories



Rio thief: How Olympic tourists can protect phone



Chicago's deadly Father's Day weekend: 12 murders. 54 shootings

the European Union, last Wednesday's edition of The Sun, the country's most popular daily newspaper, produced a revealing metaphor.

A day after the Rupert Murdoch-owned tabloid declared its support for the "Leave" campaign, its front-page went full-hysteria. A "biblical swarm" of supermoths has invaded Britain from mainland Europe, it exclaimed without a trace of irony, "and is set to annihilate our cabbages."

The prominence given to these ghastly foreign moths, accompanied on Page 2 with a helpful diagram depicting a battalion of Lepidoptera hovering over the English Channel, tells you all you need to know about Britain's EU referendum debate, where fear and frustration have come to overshadow the biggest political decision in a generation.



#### A T C J J O B C O C C d ni E 1 C o h' re d S S' P Y ( N Y Y R C 1 C 1 A S ( 1 C)

📄 edition.cnn.com/2016/06/21/opinions/brexit-campaign-disaffection-wismayer... 🎡 🔒 👼 🗢

#### International Edition + menu

Rio thief: How Olympic tourists can protect phone

Regions + CNN

Yoav

Ξ

British lawmakers clash ahead of Brexit vote 05:57

#### Story highlights

Britons will vote Thursday on whether their country should remain in the European Union

Henry Wismayer: Urge to leave has caught the wind of native disaffection

Editor's Note: Henry Wismayer is a writer and commentator based in London. You can follow him @henrywismayer. The views expressed are his own.

(CNN) — For any external observer curious to understand the issues shaping the debate over Britain's membership of



Top stories

Chicago's deadly Father's Day weekend: 12 murders. 54 shootings

the European Union, last Wednesday's edition of The Sun, the country's most popular daily newspaper, produced a revealing metaphor.

A day after the Rupert Murdoch-owned tabloid declared its support for the "Leave" campaign, its front-page went full-hysteria. A "biblical swarm" of supermoths has invaded Britain from mainland Europe, it exclaimed without a trace of irony, "and is set to annihilate our cabbages."

The prominence given to these ghastly foreign moths, accompanied on Page 2 with a helpful diagram depicting a battalion of Lepidoptera hovering over the English Channel, tells you all you need to know about Britain's EU referendum debate, where fear and frustration have come to overshadow the biggest political decision in a generation.


#### 

📄 edition.cnn.com/2016/06/21/opinions/brexit-campaign-disaffection-wismayer... ☆ 🔒 🜉

International Edition + menu

Yoav

#### CNN Regions +

British lawmakers clash ahead of Brexit vote 05:57

#### Story highlights

Britons will vote Thursday on whether their country should remain in the European Union

Henry Wismayer: Urge to leave has caught the wind of native disaffection Editor's Note: Henry Wismayer is a writer and commentator based in London. You can follow him @henrywismayer. The views expressed are his own.

(CNN) — For any external observer curious to understand the issues shaping the debate over Britain's membership of





Rio thief: How Olympic tourists can protect phone

Q.F.

Chicago's deadly Father's Day weekend: 12 murders. 54 shootings

the European Union, last Wednesday's edition of The Sun, the country's most popular daily newspaper, produced a revealing metaphor.

A day after the Rupert Murdoch-owned tabloid declared its support for the "Leave" campaign, its front-page went full-hysteria. A "biblical swarm" of supermoths has invaded Britain from mainland Europe, it exclaimed without a trace of irony, "and is set to annihilate our cabbages."

The prominence given to these ghastly foreign moths, accompanied on Page 2 with a helpful diagram depicting a battalion of Lepidoptera hovering over the English Channel, tells you all you need to know about Britain's EU referendum debate, where fear and frustration have come to overshadow the biggest political decision in a generation.

#### and maybe each paragraph is a document?

## What are we processing?

Document Collection / Cluster

Document

Section

Paragraph

Sentence

Word



#### sentences = text.split(".")?



(we'll get back to this)

## What are we processing?

Document Collection / Cluster

Document

Section

Paragraph

Sentence

Word





sequence of characters

basic unit of meaning?

white-space tokenized?



doesn't John's unlucky

#### A Turkish word

uygar\_laş\_tır\_ama\_dık\_lar\_ımız\_dan\_mış\_sınız\_casına

"as if you are among those whom we were not able to civilize (=cause to become civilized)"

uygar: *civilized* 

<u>laş</u>: become

- \_tir: cause somebody to do something
- \_ama: not able

\_dik: past participle

<u>lar</u>: plural

- \_miz: 1st person plural possessive (our)
- \_dan: among (ablative case)

\_mış: past

\_siniz: 2nd person plural (you)

\_casina: as if (forms an adverb from a verb)

K. Oflazer pc to J&M

by Julia Hockenmaier



שלהם



שלהם שלו



שלהם שלו איתם



סיפרו









he told it he<sub>1</sub> cut his<sub>2</sub> hair



מסין לסין בסין שסין כשסין ומסין ולסין



בצלם



בצלם

בצל של הם

ב צל של הם

ב צלם



בצלם

#### בצל של הם





#### 我开始写小说 = 我 开始 写 小说 *I start(ed) writing novel(s)*

by Julia Hockenmaier



#### whitespace is not enough



במהירות בעקבות



ice cream

web site

New York

New York-Based



gave up

made sense

took a picture

took apart

took the toy apart



#### Basic Structure of an NLP Algorithm





# Basic Structure of an NLP Algorithm





#### Basic Structure of an NLP Algorithm





# Basic Structure of an NLP Algorithm





# Basic Structure of an NLP Algorithm



Supervised: text 
vector
ML
Model
Model

this is not trivial!



#### Case Study: Sentence Boundary Detection.



#### sentences = text.split(".")?



A **graphic adventure game** is a form of adventure game. They are distinct from text adventures. Whereas a player must actively observe using commands such as "look" in a text-based adventure, graphic adventures revolutionized gameplay by making use of visual human perception. Eventually, the text parser interface associated with older interactive fiction games was phased out in favor of a point-and-click interface, i.e., a game where the player interacts with the game environment and objects using an on-screen cursor. In many of these games, the mouse pointer is context sensitive in that it applies different actions to different objects.



#### sentences = text.split(". ")?



Georg "Mr. George" Svendsen (19 March 1894 – 1966) was a Norwegian journalist and crime novelist.

He was born in Eidanger, and started his journalistic career in Bratsberg-Demokraten before moving on to Demokraten where he was a subeditor. In 1921 he was hired in Fremtidenand replaced in Demokraten by Evald O. Solbakken. In 1931 he was hired in Arbeiderbladet. Under the pen name "Mr. George" he became known for his humorous articles in the newspaper. At his death he was also called "the last of the three great criminal and police reporters in Oslo", together with Fridtjof Knutsen and Axel Kielland. He was also known for practising journalism as a trade in itself, and not as a part of a political career. He retired in 1964, and died in 1966.

He released the criminal novels Mannen med ljåen (1942), Ridderne av øksen (1945) and Den hvite streken (1946), and translated the book S.S. Murder by Quentin Patrick as Mord ombord in 1945. He released several historical books: Rørleggernes Fagforenings historie gjennem 50 år (1934), Telefonmontørenes Forening Oslo, gjennem 50 år (1939), Norsk nærings- og nydelsesmiddelarbeiderforbund: 25-års beretning (1948), De tause vitner: av rettskjemiker Ch. Bruffs memoarer (1949, with Fridtjof Knudsen) and Elektriske montørers fagforening gjennom 50 år (1949).



#### Sentence Boundary Detection

• How would you solve this?

(discuss)



#### Sentence Boundary Detection

• Worth reading:

Sentence Boundary Detection and the Problem with the U.S.

**Dan Gillick** Computer Science Division University of California, Berkeley dgillick@cs.berkeley.edu


## What is a sentence?

### It was a bad time. (it is always a bad time.)



# What is a sentence?

In a quiet voice, he said "this will not work. I am quitting", and the he left the room.





How often do you find yourself on the internet looking at the same boring pages? You know there is something out there but you don't know where to look

Trust me, how bad could it be?



**Disclaimer**: The websites you are forwarded to may contain naughty stuff. here is no known "porn" but sites change hands often and sometimes they change ethical direction. If you see something that looks like porn or is otherwise unacceptable let us know and we'll decide if it should be removed or if you're just a crank who likes to complain.

Don't leave the house with your fly open. Make sure you've got up-to-date virus protection.

#### Colophon & Warning & Contact

Want to add your own links? After reading the warning type the url in the box below

add

If you see something that shouldn't be here let us know

random website dot com 2001-2011

"I'd like to quit thinking of the present, like right now, as some minor insignificant preamble to something else."



## Sentence Boundary Detection

- Perhaps the most basic task.
- Non-trivial.
- Need to consider corpus, features, annotation procedure, biases.
- Need to think of your use-cases, choose your sentence definition, methods, and trade-offs.



# That's it for today.