# NATURAL LANGUAGE PROCESSING



ADOPTED SOME MATERIALS DEVELOPED IN PREVIOUS COURSES BY MANCY MCCRACKEN, LIZ LIDDY AND OTHERS; AND SOME INSTRUCTOR RESOURCES FOR THE BOOK "SPEECH AND LANGUAGE PROCESSING" BY DANIEL JURAFSKY AND JAMES H. MARTIN

#### NATURAL LANGUAGE PROCESSING (NLP)

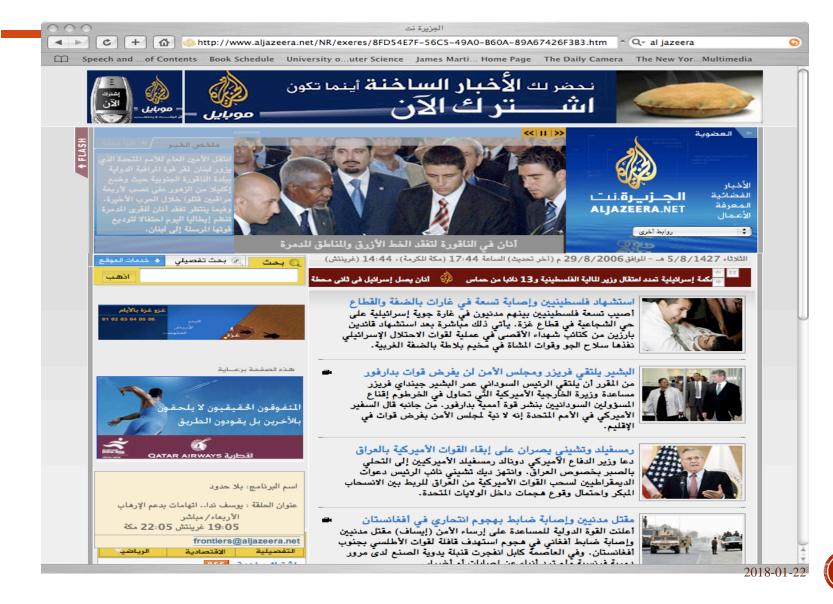
- A range of computational techniques:
  - for analyzing and representing naturally occurring texts at one or more levels of linguistic analysis
  - for the purpose of achieving human-like language processing
  - for a range of particular tasks or applications.
- Computational Linguistics doing linguistics on computers
  - Closely related, often treated as synonymous with NLP

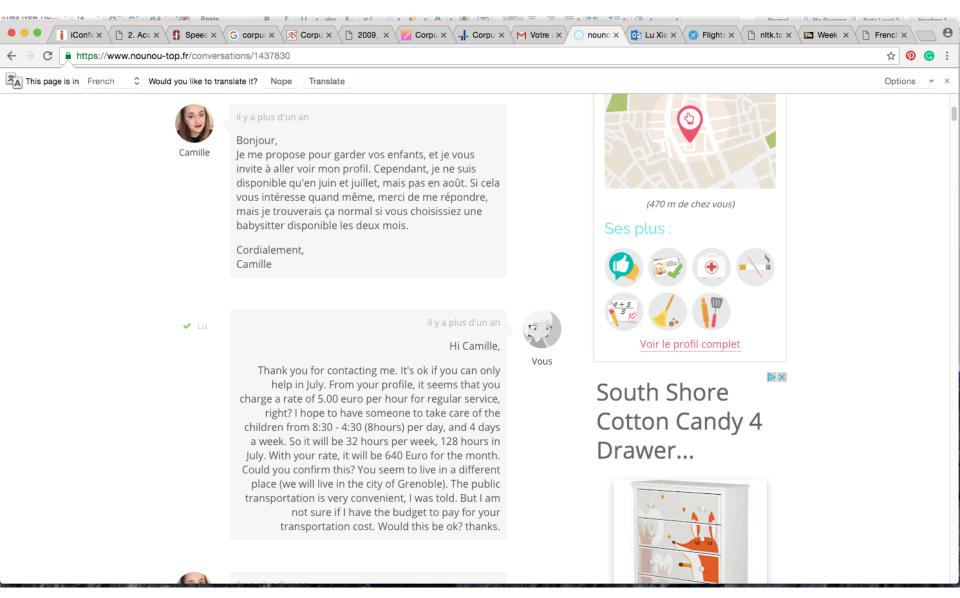
## WHERE IS NLP NOW?

- Goals can be far-reaching
  - True text understanding
  - Reasoning about knowledge in text
  - Real-time participation in spoken dialogs
- Or very down-to-earth
  - Finding the price of products on the web
  - Context-sensitive spell-checking
  - Analyzing authorship or opinions statistically
  - Extracting facts or relations from documents
  - Remembering previous searches and contexts to guide future interactions
- Currently, NLP is providing these practical applications (yet still dreaming of the AI goals)

- Information Retrieval / Search Engines provision of documents containing requested information
  - Google, many other search engines
  - Use lowest levels of NLP to stem words, find phrases for indexing documents
  - Users conform to keyword query restriction, instead of natural language queries
- Machine Translation conversion of text from one language to another
  - Usefulness of Parallel Corpora
  - Often statistically based patterns of word usage and context
  - Google, Yahoo and Bing all have language translators
  - MT techniques use context, not just word for word substitutio

## GOOGLE TRANSLATE





- Information Extraction / Text-mining populating a structured database with specific bits of information found in text
  - Competitive Intelligence analyzes news text and web blogs for
    - Names of people, companies and other entities
    - Relations between them, e.g. corporate roles, or events such as mergers

#### **Weblog Analytics**

Data-mining of Weblogs, discussion forums, message boards, user groups, and other forms of user generated media

- Product marketing information
- Political opinion tracking
- Social network analysis
- Buzz analysis (what's hot, what topics are people talking about right now).

Human-computer Interfaces – NLP assistants, chatbots, interactive

querying of databases



- <u>Summarization</u> abstraction and condensation of text's major points
  - Current systems select a set of significant sentences from the document as a summary
  - Example summarizer: <a href="http://textsummarization.net/text-2018-01-2">http://textsummarization.net/text-2018-01-2</a>



Question & Answering Systems – focused information provision



 Metadata Generation – assignment of values for metadata elements in a particular standard

Elizabeth D. Liddy, Eileen Allen, Sarah Harwell, Susan Corieri, Ozgur Yilmazel, N. Ercan Ozgencil, Anne Diekema, Nancy McCracken, Joanne Silverstein, and Stuart Sutton. 2002. Automatic metadata generation & evaluation. In *Proceedings of the 25th annual international ACM SIGIR conference on Research and development in information retrieval* (SIGIR '02). ACM, New York, NY, USA, 401-402. DOI=http://dx.doi.org/10.1145/564376.564464

### THE TRENDS

- 1. An enormous amount of knowledge is now available in machine readable form as natural language text
- 2. Conversational agents are becoming an important form of human-computer communication
- 3. Much of human-human communication is now mediated by computers

# NEED FOR (MORE) NLP WORK

- Huge amounts of data
  - Internet
  - Intranet
- Applications for processing large amounts of texts require NLP expertise
- Data Science/Text Mining

Classify text into categories
Index and search large texts
Automatic translation of web
documents in different languages
Speech understanding
Understand phone conversations
Information extraction
Extract useful information from resumes
Automatic summarization

Condense 1 book into 1 page
Daily news summaries
Question answering
Knowledge acquisition
Text generations / dialogues







least

CMV: A Trump vs Hillary election would be the worst choice we've had in a preside least)

13 hours ago by SayNoToStim

So as of today, it looks like the presidential race is shaping up to be Donald Trump vs Hillary Clinton. If thi worst combination of choices we've had in 75 years, if not more. The second worst, in my opinion, was Bu parties had some sort of merit behind them (at the time)

In my eyes Hillary is the most corrupt presidential candidate we've seen since Nixon, and even then, takin consideration, the general public wasn't aware of how c

I would start listing why Trump is a bad idea but I don't r

I would start listing why Trump is a bad idea but I don't

Both parties seem wildly disconnected to the general put technology, and quite frankly scare me when it comes to

Hello, users of CMV! This is a footnote from your modernmember to read through our rules. If you see a conspeaking of which, downvotes don't change views through our popular topics wiki first. Any questions

CMV: I think it's a selfish motive to purposely try to have children.

1 year ago by Caltybeck

I want children one day. Part of me wants to have my own children but I can't justify birthing my own kids when it's such an extremely selfish motive. Sure, once you've had the kids it's selfless because of how much you have to give up for them. But the initial desire to birth them in the first place is selfish. I want my own children because I want to carry on my own genes. I want to have a little human that resembles me. As a woman, I want to experience the feeling of a baby inside of me. These reasons are the main reasons why people choose to birth their own kids. Here's a few reasons wity I find it selfish.

- There are plenty of children out there without parents. Over 150 million orphans, not to mention foster children

- Then here are people who are trying to birth their own when there are so many helpless children without a loving home.

There are over 800 million people starving in the world. People are dying from hunger and you're trying to bring another mouth that needs feeding into the world. [+184 words]

Omega037 74∆ 7 points 1 year ago

Economically and socially, not having children when you have the means to raise them is a far more selfish act.

First, let's focus the discussion on the US, since that is where I assume you live. After all, most of those 150 million orphans cannot be easily or legally adopted into the US, and those 800 million starving people is not because we lack food (we have a major surplus), it is about losistics and failed oolitical states.

Anyways, in the US, there were only 101,666 children legally up for adoption in 2012, and of them 52,039 children were adopted. Source
That contrasts the ~3.9 million babies born each year. Source

In other words, the majority children in the US who can be adopted are adopted, and even if all of them were, it would not cover even 2.5% of the births that happen.

It is also important to note that even with all those births, the US doesn't meet its replacement rate. That means that our population would be declining if not for massive amount of immigrants we take in. [+168 words]

Persuasive power of online comments in Reddit's "Change My View" discussions

of before and that's how

e adopted? I did a llion children are with:

> ng those two things. boking from all of us.

in the womb or drug use

341 comments share

top 200 comments show all 341

sorted by: best ▼

one even if people in chargeof this kid tried their best

Are you able to love a child as your own whose basic personality might be very different to yours? (not all, but some of personality is, as far as we can tell at the moment, heritable).

Because, if you are not sure about either of this, doing the "selfish" thing might be doing the better thing. A child you adopt deserves your unconditional love as parent, just as a birth-child would.

also, why are you responsible for other people you don't know? why do you think being selfish is automatically bad?

## WHY IS NLP SO HARD?

#### Seems pretty simple for humans

 Usually quite unaware of the complexity of the language tasks they perform so effortlessly

#### Some reasons are

- Subleties of meaning
  - Irony, sarcasm, humor, metaphor
- Ambiguity
  - Ambiguity is a fundamental problem of computational linguistics
  - Resolving ambiguity is a crucial goal

## **AMBIGUITY**

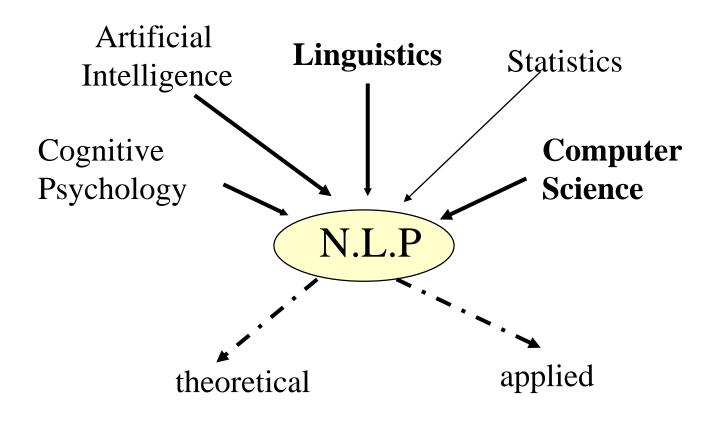
- Find at least 5 meanings of this sentence:
  - I made her duck
  - I cooked waterfowl for her benefit (to eat)
  - I cooked waterfowl belonging to her
  - I created the (plaster?) duck she owns
  - I caused her to quickly lower her head or body
  - I waved my magic wand and turned her into undifferentiated waterfowl

#### 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"

### AMBIGUITY IS PERVASIVE

- Grammar: Make can be:
  - Transitive: (verb has a noun direct object)
    - I cooked [waterfowl belonging to her]
  - Ditransitive: (verb has 2 noun objects)
    - I made [her] (into) [undifferentiated waterfowl]
  - Action-transitive (verb has a direct object and another verb)
  - I caused [her] [to move her body]



Natural Language Processing

Language Analysis\* Language Generation



#### CATEGORIES OF KNOWLEDGE

- Phonology
- Morphology
- Syntax
- Semantics
- Pragmatics
- Discourse

## COURSE OUTLINE

- Lectures
  - Mar. 26 away for a conference trip; Ph.D. student Jenna Kim will deliver the lecture and hold the lab session
- Jenna will also hold office hours in the two weeks before the assignment due date to help you on the homework
  - Wednesdays: 1:00 5:00 pm
  - Thursdays: 10:00 am 5:00 pm

# METHODS OF ASSESSMENT

Assessment method	Due Dates (submission will be due by 11:59 pm. of the day unless otherwise noted)
Class	Notes: this includes participation in weekly lecture and
Participation	Python exercises, and contributions to class discussions
Homework	l <sup>st</sup> : Feb. 18
assignments	2 <sup>nd:</sup> Mar. 18
	3 <sup>rd</sup> : Apr. 1
	4 <sup>th</sup> : Apr. 22
NLP	Group report: Apr. 29
Application	Group poster: Apr. 29
Investigations	

# NATURAL LANGUAGE TOOLKIT (NLTK) AND PYTHON



#### PYTHON AND NLP

- Python is freely available for many platforms from the Python Software Foundation:
  - http://www.python.org/
  - Named for the group Monty Python
  - We are using Python version 3.x
    - (not backward compatible with Python 2.x)





### CHARACTERISTICS OF PYTHON

- Easy-to-learn scripting language, similar in many aspects to Perl, but with WYSIWYG block structure
- Object-oriented, with modules, classes, exceptions, highlevel dynamic data types, similar to Java
- Strongly typed, but without type declarations (dynamic typing)
- Regular Expressions and other string processing features
- Many libraries offer wide functionality:
  - https://xkcd.com/353/

[optional reading about strongly typed and dynamic typed: <a href="http://stackoverflow.com/questions/11328920/is-python-strongly-typed">http://stackoverflow.com/questions/11328920/is-python-strongly-typed</a>

https://wiki.python.org/moin/Why%20is%20Python%20a%20dynamic%20language%20and%20also%20a%20strongly%20typed%20language

#### GETTING STARTED IN PYTHON

- Python can be run as an interactive system
  - Type in expressions or small pieces of programs to try them out
- or as a command-line system.
  - Run stored python programs
- For both, it is recommended to use a Python development environment
  - IDLE is standard but really simple: especially good to edit Python programs in IDLE to keep track of the indentation for block structure
    - Or try Wing free version, PyCharm or iPython to get an IDE

# NATURAL LANGUAGE TOOLKIT (NLTK)

- A suite of Python libraries for symbolic and statistical natural language programming
  - Developed at the University of Pennsylvania
- Developed to be a teaching tool and a platform for research NLP prototypes
  - Data types are packaged as classes
  - Goal of code is to be clear, rather than fastest performance
    - But increasingly production level software is made available through wrappers
- Latest version is compatible with Python 3.x
- Online book:
  <a href="http://www.nltk.org/book/">http://www.nltk.org/book/</a>
- Authors:
   Edward Loper, Ewan Kline
   and Steven Bird



#### USING NLTK IN NLP

- NL ToolKit provides libraries of many of the common NLP processes at various language levels
  - Leverage these libraries to process text
- Goal is to learn about and understand how NLP can be used to process text without programming all processes
  - However, some programming is required to
    - Call libraries
    - Process data
    - Customize NLP processes
  - Programming language is Python

### INTRODUCTION TO NLTK

- NLTK provides:
  - Basic classes for representing data relevant to Natural Language Processing.
  - Standard interfaces for performing NLP tasks such as tokenization, tagging and parsing
  - Standard implementation of each task which can be combined to solve complex problems

#### SOME NLTK MODULES

- corpora: a package containing modules of example text
- tokenize: functions to separate text strings
- probability: for modeling frequency distributions and probabilistic systems
- stem package of functions to stem words of text
- wordnet interface to the WordNet lexical resource
- chunk identify short non-nested phrases in text
- etree: for hierarchical structure over text
- tag: tagging each word with part-of-speech, sense, etc.
- parse: building trees over text
  - recursive descent, shift-reduce, probabilistic, etc.
- cluster: clustering algorithms
- draw: visualize NLP structures and processes
- contrib: various pieces of software from outside contributors



#### TUTORIALS FOR PYTHON AND NLTK

Python

many language constructs best documented in Python 2.x: https://docs.python.org/2/

Python 3.x language reference, particularly for Unicode and string representations: <a href="https://docs.python.org/3/">https://docs.python.org/3/</a>

NLTK is a SourceForge project at: <a href="http://www.nltk.org">http://www.nltk.org</a>

documentation: <a href="http://www.nltk.org/documentation">http://www.nltk.org/documentation</a>, including

book: http://www.nltk.org/book/

API: http://www.nltk.org/api/nltk.html