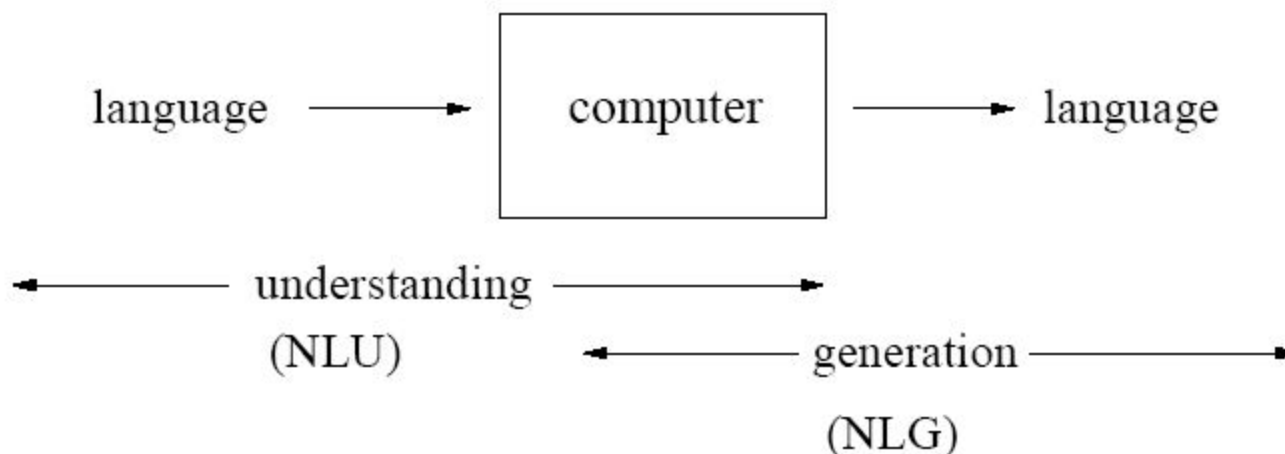


Natural language processing (NLP)

Python for AI

What is NLP?

- Natural Language Processing (NLP) is a field in Artificial Intelligence (AI) devoted to creating computers that use natural language as input and/or output.



Why NLP?



Modern NLP:

NLP in the Days of Big Data

Three trends:

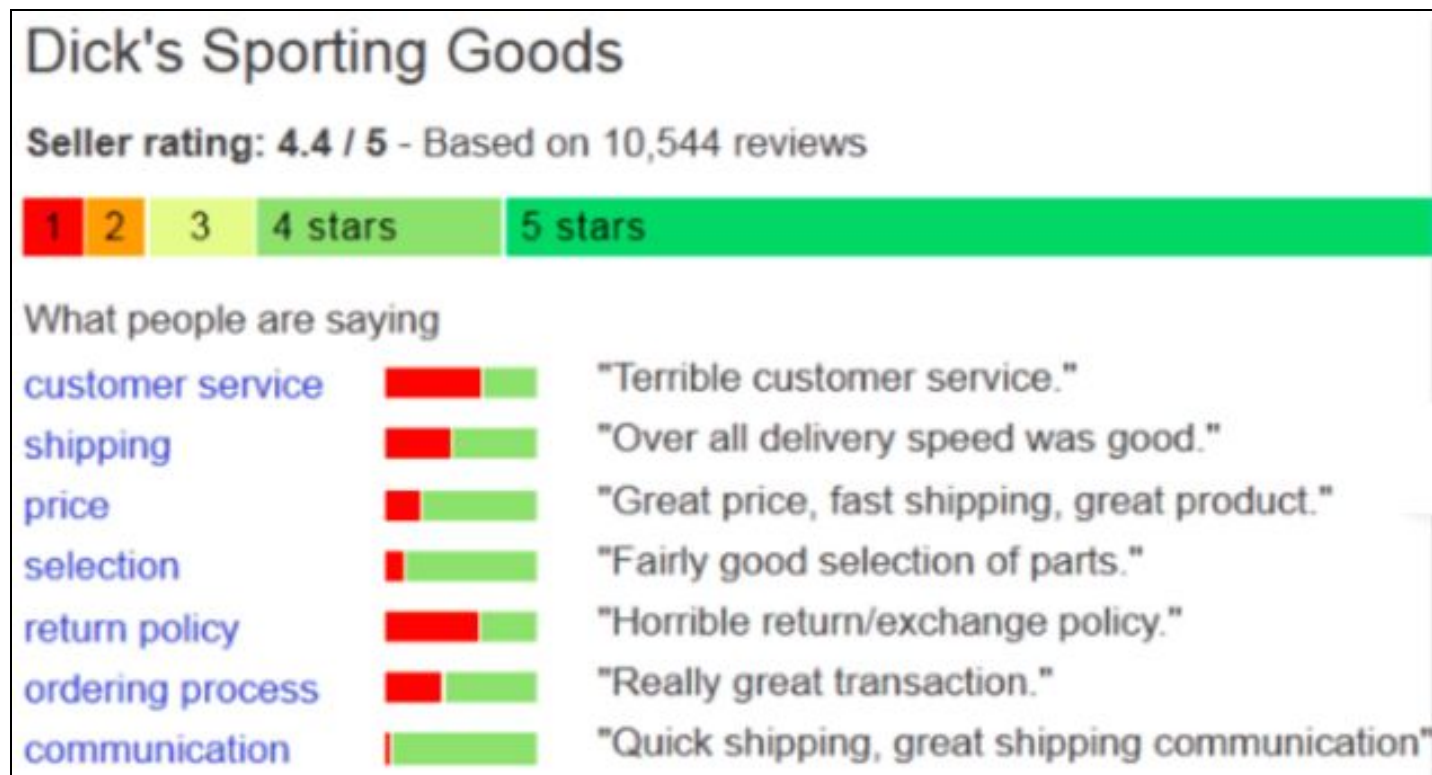
1. An **enormous amount of information** is now available in machine readable form as natural language text (newspapers, web pages, medical records, financial filings, product reviews, discussion forums, etc.)
2. Conversational agents are becoming an important form of human-computer **communication**
3. Much of human-human interaction is now mediated by computers via **social media**

NLP Applications

- Three prominent application areas:
 - Text analytics/mining (from “***unstructured data***”)
 - Sentiment analysis
 - Topic identification
 - Digital Humanities (“*new ways of doing scholarship that involve collaborative, transdisciplinary, and computationally engaged research, teaching, and publishing.*”)
 - Conversational agents
 - Siri, Cortana, Amazon Alexa, Google Assistant
 - Chatbots
 - Machine translation

Text Analytics

- Data-mining of weblogs, microblogs, discussion forums, user reviews, and other forms of user-generated media.



Text Analytics (cont.)

- Typically this involves the extraction of **limited** kinds of semantic and pragmatic information from texts
 - Entity mentions
 - Concept identification
 - Sentiment

The screenshot displays the 'API TEST TOOL' interface. At the top, there are three dropdown menus: 'English', 'Entities', and 'Graphical'. Below these, a text box contains a sample paragraph with various entities highlighted in colored boxes. To the right of the text box is a 'LEGEND color key' titled 'ENTITIES' which lists 10 types of entities with corresponding color-coded icons: Person name (pink), Car license plate (magenta), Place (purple), Phone number (orange), Email address (brown), Date (gold), Hour (grey), Money (dark grey), Address (dark purple), and Twitter hashtag (light brown).

English Entities Graphical

I really enjoyed using the **Canon Ixus** in **Madrid** on **March 4**. The **Panasonic Lumix** is a bit disappointing, but the **Canon** camera is not bad at all. All I want when taking photos is point it and then just press the button. For only **200 dollars**, a really fair price, this camera is perfect for me. Besides, I have had a good customer service experience. **John Faraday** was very nice!

LEGEND color key ENTITIES

Type of entities:

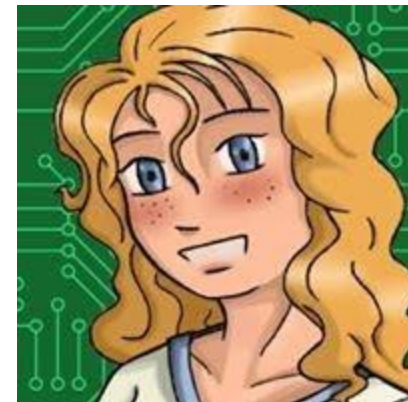
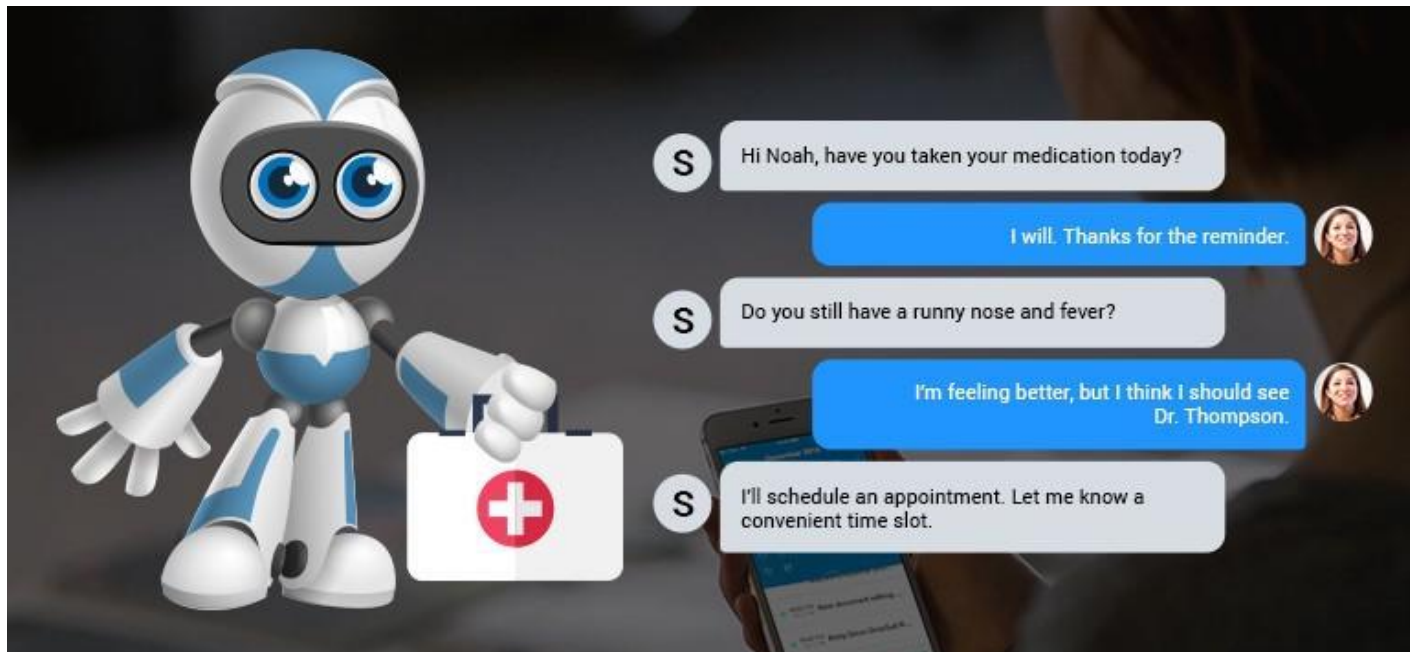
Person name	Date
Car license plate	Hour
Place	Money
Phone number	Address
Email address	Twitter hashtag

Demo

- Sentiment Analysis with Python NLTK Text Classification
 - <http://text-processing.com/demo/sentiment/>
- Tweet Sentiment Visualization Tool
 - https://www.csc2.ncsu.edu/faculty/healey/tweet_viz/tweet_app/
- Concept Extraction
 - <http://aylien.com/concept-extraction/>

Conversational Agents

- Combine
 - Speech recognition/synthesis
 - Question answering
 - From the web and from structured information sources (freebase, dbpedia, yago, etc.)
 - Simple agent-like abilities
 - Create/edit calendar entries
 - Reminders
 - Directions
 - Invoking/interacting with other apps



Mitsuku

Question Answering

- Traditional *information retrieval* provides documents/resources that provide users with what they need to satisfy their information needs.
- *Question answering* on the other hand directly provides an answer to information needs posed as questions.

IBM Watson



https://www.youtube.com/watch?v=WFR3lOm_xhE

Machine Translation

- The automatic translation of texts between languages is one of the oldest non-numerical applications in Computer Science.
- In the past 15 years or so, MT has gone from a niche academic curiosity to a robust commercial industry.

巨大な銃規制集会が米国を席卷

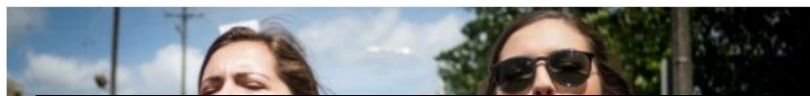
学生が主催する「私たちの生活のための行進」イベントでは、全国的に数十万人の抗議者が集まります。

🕒 4時間 | 米国とカナダ

Huge gun-control rallies sweep US

Student-led March For Our Lives events nationwide draw hundreds of thousands of protesters.

🕒 4h | US & Canada



But NLP very is hard..

- Understanding natural languages is hard ... because of inherent *ambiguity*
- Engineering NLP systems is also hard ... because of:
 - Huge amount of data resources needed (e.g. grammar, dictionary, documents to extract statistics from)
 - Computational complexity (intractable) of analyzing a sentence

Ambiguity (1)

“Get the cat with the gloves.”



Ambiguity (2)

Find at least 5 meanings of this sentence:

“I made her duck”

1. I cooked waterfowl for her benefit (to eat)
2. I cooked waterfowl belonging to her
3. I created the (plaster?) duck she owns
4. I caused her to quickly lower her head or body
5. I waved my magic wand and turned her into undifferentiated waterfowl

Ambiguity is Pervasive

- **Phonetics**

- I mate or duck
- I'm eight or duck
- Eye maid; her duck
- Aye mate, her duck
- I maid her duck
- I'm aid her duck
- I mate her duck
- I'm ate her duck
- I'm ate or duck
- I mate or duck

Sound like

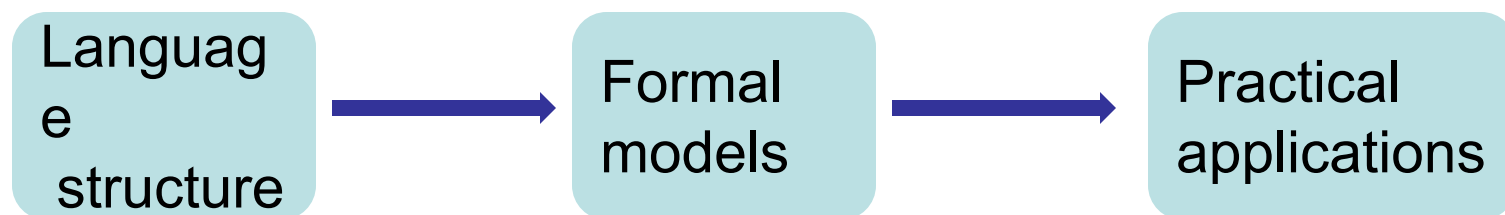
"I made her duck"

The Bottom Line

- Complete NL Understanding (thus general intelligence) is impossible.
- But we can make incremental progress.
- Also we have made successes in **limited domains**.

The Big Picture Approach

All of these applications operate by **exploiting underlying regularities** in human languages. Sometimes in complex ways, sometimes in pretty trivial ways.



Topics: Linguistics

- Word-level processing
- Syntactic processing
- Lexical and compositional semantics
- Discourse structure