

Speech and Language Processing

An introduction to the Natural Language Processing course

Ing. R. Tedesco. PhD, AA 20-21

(mostly from: Speech and Language Processing - Jurafsky and Martin)

Why Should You Care?

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

Major Topics

1. Words

2. Syntax

3. Meaning

4. Discourse

5. Speech

6. Applications exploiting each




Big Applications

- Question answering
- Conversational agents
- Summarization
- Machine translation

Topics: Linguistics

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

Topics: Techniques

- Finite-state methods
 - Context-free methods
 - Augmented grammars
 - ◆ Unification
 - ◆ Lambda calculus
 - First order logic
- 
- Probability models
 - Supervised machine learning methods
 - Neural Networks

Topics: Applications

- Small
 - ◆ Spelling correction
 - ◆ Hyphenation
- Medium
 - ◆ Word-sense disambiguation
 - ◆ Named entity recognition
 - ◆ Information retrieval
- Large
 - ◆ Question answering
 - ◆ Conversational agents
 - ◆ Machine translation
- Stand-alone
- Enabling applications
- Funding/Business plans

Categories of Knowledge

- Phonology
 - Morphology
 - Syntax
 - Semantics
 - Pragmatics
 - Discourse
 - Prosody
- Each kind of knowledge has associated with it an encapsulated set of processes that make use of it.
- Interfaces are defined that allow the various levels to communicate.
- This usually leads to a pipeline architecture.

Ambiguity

- Computational linguists are obsessed with 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]

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

Brief history - 1

Starting research fields:

- ◆ Linguistics
- ◆ Natural Language Processing (computer science)
- ◆ Speech Recognition (electronics)
- ◆ Computational Linguistics (psychology)

1940-1950 - World War II

- ◆ Finite State Automata: Formal Language Theory (algebra and set theory for the formalization of languages) - Chomsky, Backus and Naurs
- ◆ Probabilistic algorithms for speech, information theory (Shannon), noise of the channel encoding and decoding, entropy of a language
- ◆ Machine Translation is the most desired application

Brief history - 2

1957-1970 - Two paradigms

◆ **Symbolic**

- a) Formal Language Theory (Chomsky): parsing algorithms (first top-down and bottom-up then with dynamic programming)
- b) Artificial Intelligence - Logic Theories - (from Newell and Simon) - a combination of pattern matching and search for keywords with simple heuristics for reasoning and answer questions
- a) and b) lead to the early systems

- ◆ **Stochastic:** Bayesian method and use of dictionaries and corpora (the first OCR) - Browning, Mosteller and Wallace. The Brown Corpus - Kucera and Francis

Brief history - 3

1970-1983 – FS Models

- ◆ Understanding natural language - Winograd (the SHRDLU parser and construction of a systemic grammar): parsing well understood
- ◆ You could start working seriously on semantics and discourse (Schank et al: scripts, plans and goals, human memory) (Quillian, Rumelhart and Norman, Simmons, ...) with network-based semantics integrated 'case roles' (Fillmore)
- ◆ Discourse Modeling
 - Analysis of substructures of discourse (Grosz, Sidner)
 - Automatic resolution of references (Hobbs)
 - Belief-Desire-Intention (Perrault, Allen - Cohen and Perrault)

Brief history - 4

1983-1993 - empiricism and FS Models

- ◆ Continuation of Finite-State models
 - For phonology and morphology (Kaplan and Kay)
 - For syntax (Church)
- ◆ Return to empiricism
 - Work at IBM for speech recognition based on probabilistic models
 - Data-driven approaches: POS tagging, parsing and annotation, for ambiguity resolution, use of connectionist models... from speech recognition to the semantics
- ◆ Natural Language Generation

Brief history - 5

1994-1999 - decline of symbolic approach

- ◆ Difficulties with symbolic approach to improve
- ◆ Heavy use of data-driven methods and probabilistic models
- ◆ Enlargement of the application fields (from the Web to Alternative and Augmentative Communication...)

2000-2010 - empiricism and Machine Learning

- ◆ The empirical approach becomes even more significant
 - A lot of material written and talked about a lot of material available online and already 'annotated' (in terms of syntactic, semantic and pragmatic aspects)
- ◆ Close liaison with the research community of 'machine learning'
 - Focus on learning
 - New opportunities relied to high-performance computing resources
 - Unsupervised systems become more important than the first favorite supervised systems (the trend is set to grow further)

Brief history - 6

2010-... - Machine Learning everywhere

- ◆ Neural Networks for NLP
 - NN-based ASR/TTS, ...
- ◆ Conversational Agents
- ◆ Emotion and Affect
- ◆ Subjectivity and Sentiment Analysis
- ◆ Personality

General info on the NLP course

| Data | Dove | 09:00 | 10:00 | 11:00 | 12:00 | 13:00 | 14:00 | 15:00 | 16:00 | 17:00 |
|-----------|--|-------|-------|-------|-------|-------|-------|-------|-------|-------|
| Lunedì | | | | | | | | | | |
| Martedì | | | | | | | | | | |
| Wednesday | 9.0.3 | | | | | | | | | |
| Thursday | AULA VIRTUALE | | | | | | | | | |
| Friday | AULA VIRTUALE 9.0.3 | | | | | | | | | |
| Sabato | | | | | | | | | | |

This very first lecture: introduction (held just once, for all students)



NATURAL LANGUAGE PROCESSING esercitazione Squadra1 (dal 16/09/2020 al 04/11/2020)

Topic "A", for all students



NATURAL LANGUAGE PROCESSING lezione (dal 17/09/2020 al 05/11/2020)

Topic "B", for all students



NATURAL LANGUAGE PROCESSING lezione (dal 18/09/2020 al 30/10/2020)

Topic "C" for group 2 (person code: even number)



NATURAL LANGUAGE PROCESSING esercitazione Squadra2 (dal 18/09/2020 al 30/10/2020)

The following week...

| | | | | | | | | | | |
|-----------|-----------------------|--|--|--|--|--|--|--|--|--|
| Wednesday | 9.0.3 | | | | | | | | | |
|-----------|-----------------------|--|--|--|--|--|--|--|--|--|

Topic "C", again, for group 1 (person code: odd number)



NATURAL LANGUAGE PROCESSING esercitazione Squadra1 (dal 16/09/2020 al 04/11/2020)

And so on...

General info on the NLP course

| Data | Dove | 09:00 | 10:00 | 11:00 | 12:00 | 13:00 | 14:00 | 15:00 | 16:00 | 17:00 |
|-----------|-------------------------------|-------|-------|-------|-------|-------|-------|-------|-------|-------|
| Lunedì | | | | | | | | | | |
| Martedì | | | | | | | | | | |
| Wednesday | 9.0.3 | | | | | | | | | |
| Thursday | AULA VIRTUALE | | | | | | | | | |
| Friday | AULA VIRTUALE | | | | | | | | | |
| | 9.0.3 | | | | | | | | | |
| Sabato | | | | | | | | | | |

| | | |
|---|---|---|
| Proposal: 14:30-16:00 (no coffee break) | → | NATURAL LANGUAGE PROCESSING esercitazione Squadra1 (dal 16/09/2020 al 04/11/2020) |
| Proposal: 10:30-13:00 (15' coffee break) | ← | NATURAL LANGUAGE PROCESSING lezione (dal 17/09/2020 al 05/11/2020) |
| Proposal: 8:30-11:00 (15' coffee break) | ← | NATURAL LANGUAGE PROCESSING lezione (dal 18/09/2020 al 30/10/2020) |
| Proposal: 14:30-16:00 (no coffee break) | → | NATURAL LANGUAGE PROCESSING esercitazione Squadra2 (dal 18/09/2020 al 30/10/2020) |

Exam

- Written exam
- 3 topics, with 3 open questions each
- Could require to solve simple numeric exercises; no calculator is needed

General info on the NLP course

- Tools:
 - ◆ Python3
 - ◆ NLTK: <https://www.nltk.org>
- Web site: <http://corsi.dei.polimi.it/nlp>
- **Usually**, slide posted **before** the lecture
- Lectures (both on-line and in classroom) will be recorded
- Link to the video recording posted after the lecture
- Lectures in the classroom will be available via a live stream