

# Natural Language Processing Introduction



August 2019

# Who am I?

Mario Alejandro Campos Soberanis (mcampos@soldai.com)



Master in Computer Science (UADY)

Research Engineer at SoldAI

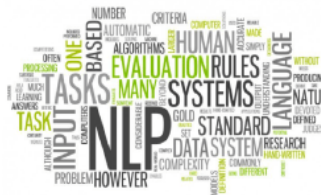
Experience as webmaster, chief technology officer and research engineer at SoldAI

Interest in conversational systems, automatic reasoning and biologically inspired algorithms

# What NLP is all about?

# Natural Language Processing

“Natural Language Processing (NLP) is the interdisciplinary field of study between artificial intelligence, linguistics and computer science whose goal is to make computers perform useful tasks that involve human language”



# What is NLP used for?

- Allow communication between human and machine (conversational agents)

# What is NLP used for?

- Allow communication between human and machine (conversational agents)
- Enhance communication between humans (automatic translation)

# What is NLP used for?

- Allow communication between human and machine (conversational agents)
- Enhance communication between humans (automatic translation)
- Make useful processing of text and speech (orthographic corrector)

# Why is relevant?

- Increment of enterprise, commercial and industrial applications using natural language (search, publicity, traslation, automatic speech recognition, automatic support).

# Why is relevant?

- Increment of enterprise, commercial and industrial applications using natural language (search, publicity, traslation, automatic speech recognition, automatic support).
- Human language as universal communication paradigm (Siri, Google Assistant, Cortana, Messenger, Alexa).



# Why is relevant?

- Increment of enterprise, commercial and industrial applications using natural language (search, publicity, traslation, automatic speech recognition, automatic support).
- Human language as universal communication paradigm (Siri, Google Assistant, Cortana, Messenger, Alexa).
- Tool to obtain knowledge of a bunch of unstructured data.

# Why is hard?

Human language is an efficient communication system, most things are not sayed but implicitly understood

# Why is hard?

Human language is an efficient communication system, most things are not said but implicitly understood

- Computational representation complexity
  - People use discourse, computers data and commands (NLP tries to close that gap)

# Why is hard?

Human language is an efficient communication system, most things are not said but implicitly understood

- Computational representation complexity
  - People use discourse, computers data and commands (NLP tries to close that gap)
- Human language is inherently ambiguous
  - Evolution, slang, regionalisms
  - Humor and sarcasm
  - Writing and grammatical errors

# Why is hard?

Human language is an efficient communication system, most things are not said but implicitly understood

- Computational representation complexity
  - People use discourse, computers data and commands (NLP tries to close that gap)
- Human language is inherently ambiguous
  - Evolution, slang, regionalisms
  - Humor and sarcasm
  - Writing and grammatical errors

The perfect understanding of the human language is an AI-complete problem.

# Turing Test

*“A computer can be considered intelligent if it’s able to hold a conversation with a human being without realizing to be talking with a machine”*

— Alan Turing

# Ambiguity

- I saw the mountains flying to New York
- After the death, the miners refuse to work
- In Mexico a woman gives birth every 15 minutes
- The officer shot the man with the knife

# Lost in translation

"The spirit is willing, but the flesh is weak"  
Translates to:



# Lost in translation

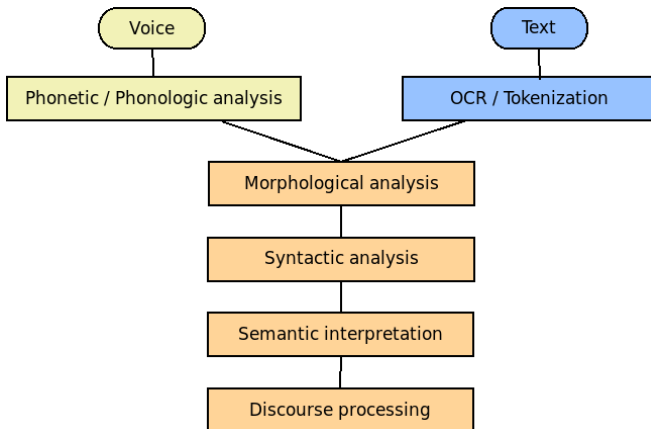
"The spirit is willing, but the flesh is weak"  
Translates to:

"The vodka is agreeable, but the meat is rotten"

# Combinatorial space for words

- A highschool student knows around 60,000 words
- Almost each sentence produced by a person is a combination generated for the first time in it's life.

# NLP Levels



# Pre-processing

- Cleaning
  - Deletion of empty meaning words (stopwords)
  - Capitalization
  - Processing of characters and symbols

# Pre-processing

- Cleaning
  - Deletion of empty meaning words (stopwords)
  - Capitalization
  - Processing of characters and symbols
- Annotation
  - Grammar labeling (Part Of Speech tagging)
  - Structural marking

# Pre-processing

- Cleaning
  - Deletion of empty meaning words (stopwords)
  - Capitalization
  - Processing of characters and symbols
- Annotation
  - Grammar labeling (Part Of Speech tagging)
  - Structural marking
- Normalization
  - Stemming
  - Lematizing

# Pre-processing

- Cleaning
  - Deletion of empty meaning words (stopwords)
  - Capitalization
  - Processing of characters and symbols
- Annotation
  - Grammar labeling (Part Of Speech tagging)
  - Structural marking
- Normalization
  - Stemming
  - Lematizing
- Others
  - Tokenizing / Segmentation
  - Counting and grouping

# Main approaches

- Rule based methods
  - Regular expressions
  - Free context Grammars
  - First order logic



# Main approaches

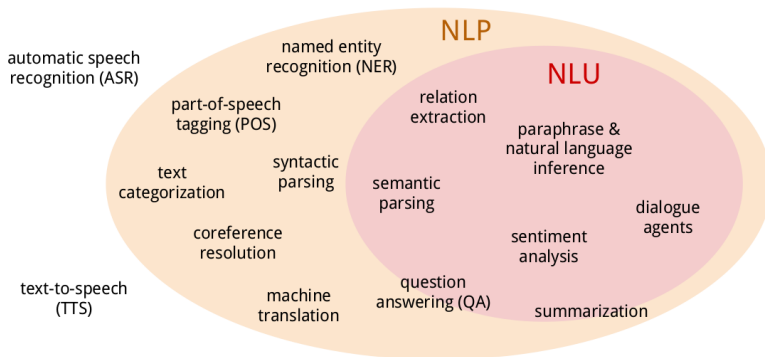
- Rule based methods
  - Regular expressions
  - Free context Grammars
  - First order logic
- Probabilistic models and Machine Learning
  - Maximum likelihood
  - Linear classification
  - Markov hidden chains

# Main approaches

- Rule based methods
  - Regular expressions
  - Free context Grammars
  - First order logic
- Probabilistic models and Machine Learning
  - Maximum likelihood
  - Linear classification
  - Markov hidden chains
- Deep learning
  - Representation Learning
  - Embeddings
  - Convolutional, Recursive, Long Short Term Memory and Recurrent Neural Networks

# Task terminology

## NLU vs. NLP vs. ASR



# Some interesting applications

- Sentiment Analysis
- Ortographic correction
- Search engines
- Information extraction
- Document classification
- Automatic translation
- Dialog systems and digital assistants
- Automatic question answering
- Natural language database interfaces
- Automatic summary

# Resources

- Libraries
  - NLTK (Natural Language Toolkit)
  - Stanford CoreNLP
  - Apache OpenNLP
  - Spacy
- Corpus and databases
  - WordNet
  - Penn TreeBank

# Let's code

Execute in your machine

```
git clone https://github.com/MaxSob/nlp-introduction
```

# References

- [1] Jurafsky, D., Martin, J.: Speech and Language Processing 2nd. ed. (2009).
- [2] Mikolov, T., Corrado, G., Chen, K. y Dean, J.: Efficient Estimation of Word Representations in Vector Space. In: CoRR, (2013).  
<https://dblp.org/rec/bib/journals/corr/abs-1301-3781>
- [3] Pinker, S.: The Stuff of Thought - Language as a window into human nature. <https://www.youtube.com/watch?v=5S1d3cNge24>