

Natural Language Processing

Unit 1: Introduction and Concepts



August 2019

Who am I?

Mario Alejandro Campos Soberanis (mario.campos@upy.edu.mx)



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

Course Syllabus

- Unit 1: Classical approaches: 11/10/2019
 - Introduction and concepts
 - Preprocessing
 - Lexical analysis
 - Syntactic analysis
 - Semantic analysis

Course Syllabus

- Unit 1: Classical approaches: 11/10/2019
 - Introduction and concepts
 - Preprocessing
 - Lexical analysis
 - Syntactic analysis
 - Semantic analysis
- Unit 2: Statistical approaches: 15/11/2019
 - Corpus
 - Classification
 - Part of speech tagging

Course Syllabus

- Unit 1: Classical approaches: 11/10/2019
 - Introduction and concepts
 - Preprocessing
 - Lexical analysis
 - Sintactic analysis
 - Semantic analysis
- Unit 2: Statistical approaches: 15/11/2019
 - Corpus
 - Classification
 - Part of speech tagging
- Unit 3: Deep learning: 13/12/2019
 - Sentiment analysis
 - Neural networks
 - Deep learning approaches

Course Syllabus

- Unit 1: Classical approaches: 11/10/2019
 - Introduction and concepts
 - Preprocessing
 - Lexical analysis
 - Sintactic analysis
 - Semantic analysis
- Unit 2: Statistical approaches: 15/11/2019
 - Corpus
 - Calssification
 - Part of speech tagging
- Unit 3: Deep learning: 13/12/2019
 - Sentiment analysis
 - Neural networks
 - Deep learning approaches

About the course

- Homework
 - Individual
 - Teams (2 persons)

About the course

- Homework
 - Individual
 - Teams (2 persons)
- Evaluation
 - Participation
 - Assignments (40 %)
 - Exam and projects (60 %)

About assignments

- Deadline weekly (if an assignment is requested on monday the deadline is next monday before 23:59:59 email/schoology time)

About assignments

- Deadline weekly (if an assignment is requested on monday the deadline is next monday before 23:59:59 email/schoology time)
- Format
 - Reports/Essays/Presentations: PDF
 - Programming assignments: Jupyter Notebook (.ipynb)
 - Projects: Python code (.py)

About assignments

- Deadline weekly (if an assignment is requested on monday the deadline is next monday before 23:59:59 email/schoology time)
- Format
 - Reports/Essays/Presentations: PDF
 - Programming assignments: Jupyter Notebook (.ipynb)
 - Projects: Python code (.py)
- Naming Individual:
NLP_{*homework_no*}-{*last_name*}-{*first_name*}.{*file_extension*}
Team:
NLP_{*homework_no*}-{*team*}-{*last_names*}.{*file_extension*}
examples: NLP_01_Campos_Mario.pdf,
NLP_03_TeamA_Campos_Soberanis_Perez.pdf



What is NLP used for?

- Allow communication between humans and machines (conversational agents)

What is NLP used for?

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

What is NLP used for?

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

Why is relevant?

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

Why is relevant?

- Increment of enterprise, commercial and industrial applications using natural language (search, publicity, translation, 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, translation, 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

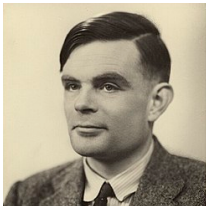
Why is hard?

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

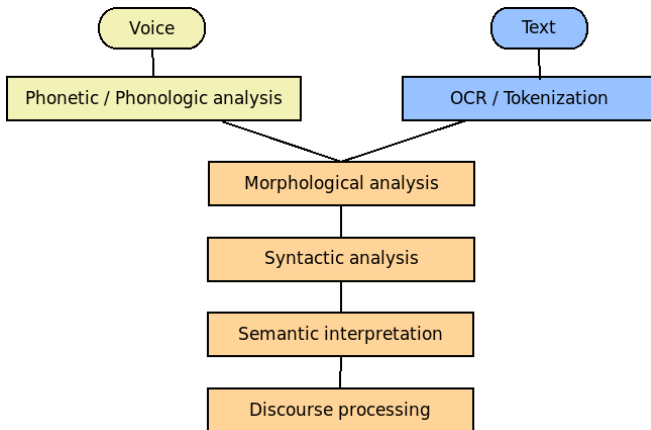
"The spirit is willing, but the flesh is weak"
Translates to:

"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

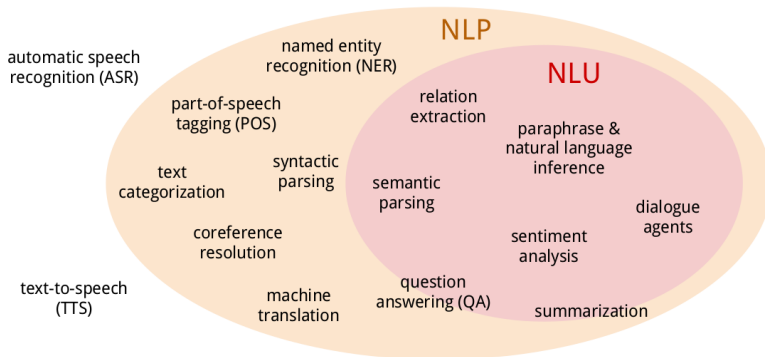
- Rule based methods
 - Regular expressions
 - Free context Grammars
 - First order logic
- Probabilistic models and Machine Learning
 - Maximum likelihood
 - Linear classification
 - Markov hidden models

Main approaches

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

Assignments

Assignment 1: Write a report about one of the following NLP tasks:

- Automatic speech recognition
- Dialogue agents
- Sentiment analysis
- Question answering

The report will include:

- Applications
- Approaches to solve the task
- Commercial products using it
- References

Let's code

```
78 // Trim the path and replace backslashes with forward slashes
79 $SESSION['captcha']['config'] = serialize($captcha_config);
80
81 return array(
82     'code' => $captcha_config['code'],
83     'image_src' => $image_src
84 );
85 }
86
87
88 if (function_exists('hex2rgb')) {
89     function hex2rgb($hex_str) {
90         $hex_str = preg_replace('/#/', '', $hex_str);
91         $rgb_array = preg_replace("/[0-9A-Fa-f]/", '', $hex_str); // Gets a proper hex string
92         if (strlen($hex_str) == 6) {
93             $color_val = hexdec($hex_str);
94             $rgb_array['r'] = 0xFF & ($color_val >> 0x10);
95             $rgb_array['g'] = 0xFF & ($color_val >> 0x8);
96             $rgb_array['b'] = 0xFF & $color_val;
97         } elseif (strlen($hex_str) == 3) {
98             $rgb_array['r'] = hexdec(str_repeat(substr($hex_str, 0, 1), 2));
99             $rgb_array['g'] = hexdec(str_repeat(substr($hex_str, 1, 1), 2));
100             $rgb_array['b'] = hexdec(str_repeat(substr($hex_str, 2, 1), 2));
101         } else {
102             return false;
103         }
104     }
105 }
106 // Draw the image
107 if (isset($_GET['code'])) {
108     $code = $_GET['code'];
109     $image_src = $captcha_config['image_src'];
110     // Draw the image
111     // ...
```


References

- [1] Jurafsky, D., Martin, J.: Speech and Language Processing 2nd. ed. (2009).
- [2] Bird S., Klein E., Loper, E.: Natural Language Processing with Python, (2009). ISBN: 978-0-596-51649-9
- [3] Indurkha N., Damerau F. Handbook of Natural Language Processing, Second Edition (2010). ISBN: 978-1-4200-8593-8
- [4] Kao, A., Poteet S. Natural Language Processing and Text Mining (2007). ISBN: 78-1-84628-175-4
- [5] 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>
- [6] Pinker, S.: The Stuff of Thought - Language as a window into human nature. <https://www.youtube.com/watch?v=5S1d3cNge24>