# Machine Learning for Natural Language Processing

## The *Why* and *What* of NLP

**Lecture 1**
**February 5th 2020, ENSAE Paris**

**Benjamin Muller**

INRIA Paris - ALMANACH
benjamin.muller@inria.fr

## This course

- We will cover techniques used in industry (Facebook, Google, Apple, Twitter...)

- Introduce core ideas at the basis of modern NLP algorithms

- Focus on machine learning applied to NLP

**Goal**: Provide a toolkit of concepts and methods to describe and tackle NLP problems in real-life.

# Course Logistics

- 6 sessions

- 1h30 lecture followed by 1h30 applied *lab* session

- Course Material nlp-ensae.github.io

# Course Outline

1. The Why and What of Natural Language Processing
2. Representing text with vectors
3. Task specific Modeling of Text
4. Neural Natural Language Processing
5. Language Modeling
6. Transfer Learning with Neural Modeling for NLP

# Course Evaluation

- Project: Implement NLP algorithm (list of projects given later)

- Outcome : Self-contained **notebook** uploaded to **github** or **google colab**

# Today session outline

- Why is language hard ? the 4 challenges of NLP
- What is Natural Language Processing ?
  - A non-exhaustive definition of NLP
  - A brief history of NLP
  - NLP in three pipelines

# Why Natural Language Processing ?

# Survival Guide

- Always asks *why ?*
- Be focused: Focus means being active (ask questions, take notes, ...)
- Practice (code) often

# Why Natural Language Processing ?

What do we do with language ?

- We communicate using language
- We think (mostly) with language
- We tell stories in language
- We describe our theories in language

Why NLP ?

- Information Retrieval (search, recommendation, aggregation)
- Better interfaces (human-computer, human-human interface)
- Better understanding of our thinking process and of language itself

# Why Natural Language Processing ?

Amount of online textual data...[1]

- 60 billion web-pages online (1.7 billion websites)
- 48,731,540 Wikipedia pages (open source)

...growing at a fast pace

- 8000 tweets/second
- 2.8 million mail / second (60% spam)
- +500 users / second

---

[1]internet live stats

# Why Natural Language Processing ?

Potential Users of Natural Language Processing

- 7.7 billion people use some sort of language (January 2019)
- 4.4 billion people connected (January 2019)

# Why Natural Language Processing ?

What products ?

- Search: $+2$ billion people use Google, 700 millions people use Baidu
- Social Media: $+3$ billion users of Social media (Facebook, instagram, WeChat, Twitter...)
- Voice assistant: $+100$ million users (Alexa, Cortona, Siri, Google Assistant)

# Why Natural Language Processing ?

Myth or Reality of ""Artificial General Intelligence"" ?

- Billions $ invested in research in AI
- Fast adoption paced : Incremental progress in research is quickly spreading to users
- Myth or Reality of AGI ?

# Why Natural Language Processing ?



Figure: The Hype Cycle

# Objective of the course

- Toolkit for how to approach any NLP problem
- Get a theoretical understanding of most recent NLP models
- Grasp the challenges (model, data, computation, time...) of NLP projects

# Why is language hard ?

# A Definition of Language

Definition 1: *Language is a means to communicate, it is a* **semiotic** *system. By that we simply mean that it is a set of signs. A sign is a pair consisting in [...] a signifier (or exponent) and a signified (or meaning).*

Definition 2: *A sign consists in a* **phonological** *structure, a* **morphological** *structure, a* **syntactic** *structure and a* **semantic** *structure*[2]

---

[2](Kracht)

# Quick introduction to linguistics

# Quick introduction to linguistics

6 Levels of analysis

- Phonological level
- Morphological Level
- Syntactic level
- Semantic Level
- Linguistic Context
- Extra-linguistic level

$\rightarrow$ All NLP problems can be split between one or several of these level of analysis

# Why is language hard ?

- Language **diversity**
- Language **variation**
- Language **ambiguity**
- Language **sparsity**

# Phonological Diversity

- Syllables are formed of phoneme sequences
- In most languages, some syllables are valid, some are not

E.g : Japanese has only one *liquid* phoneme /r/
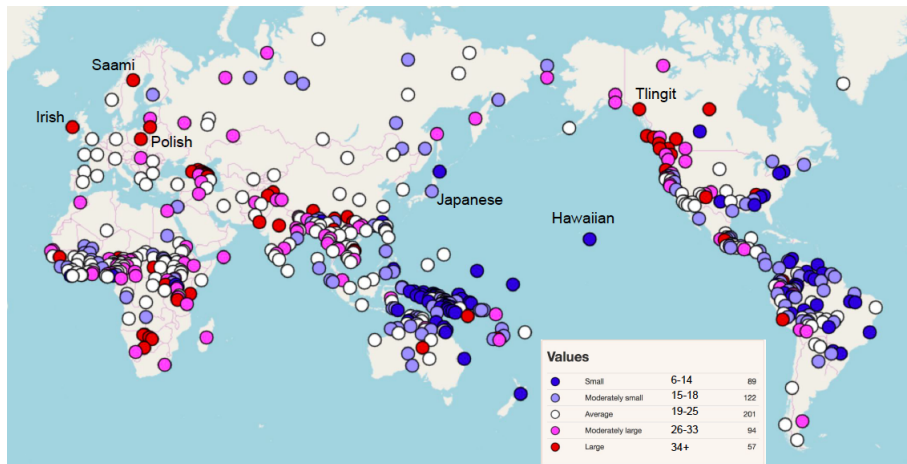
# Phonological Diversity



Figure: Consonant Inventory (size of the set of consonants) Source: The World Atlas of Language Structures

# Morphological Diversity

- Analytic and isolating languages
  - Each word carries exactly one meaning
  - e.g Chinese
- Synthetic languages
  - Agglutinative
    - Each word can have several morphs, each carrying one meaning
    - e.g : Turkish el-ler-imiz-in (HAND-pl-poss1pl-genitive) 'of our hands'
  - Fusional : - Each word can have several morphs, each carrying one or more meanings, of which (generally) only one lexical morph
  - Polysynthetic - Each word can have several lexical or grammatical morphs
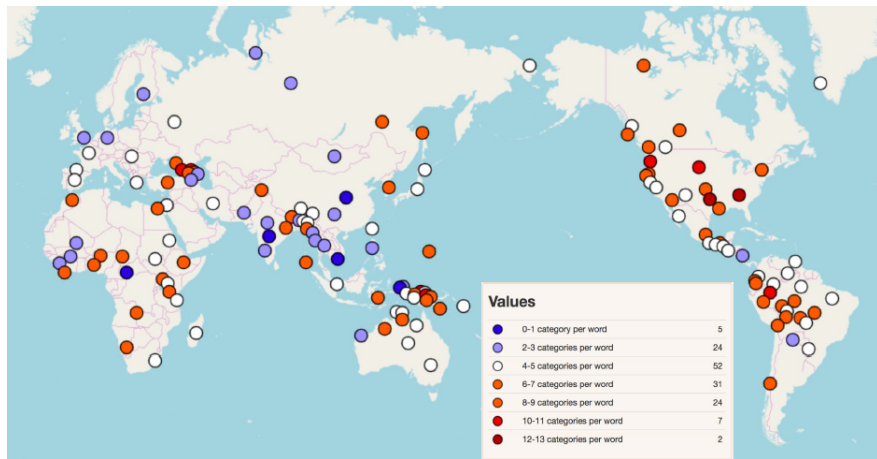
# Morphological Diversity



Figure: Number of Category per Word
Source: ►The World Atlas of Language Structures

# Syntactic Diversity

- Word order differs across languages
- Word order degree of freedom also differs across languages
- We characterize word orders with : Subject - Verb - Object order
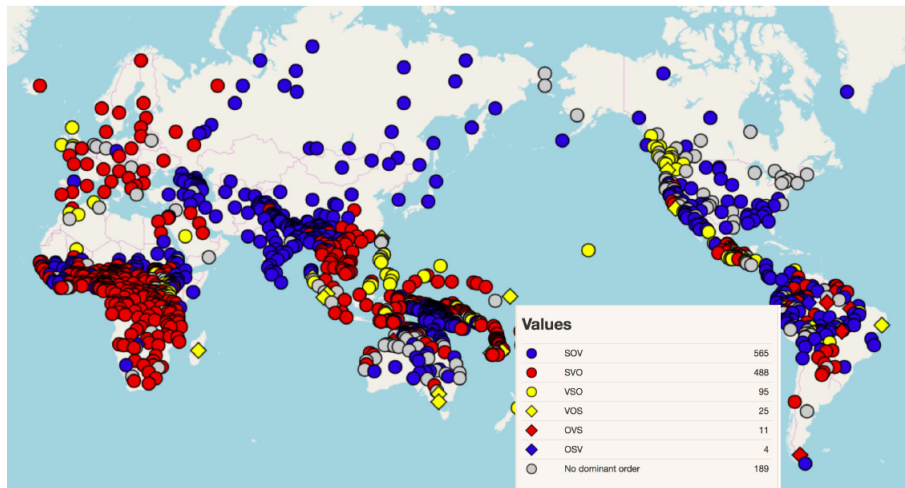
# Syntactic Diversity



Figure: Order of Subject (S) Object (O) and Verb (V)
Source: The World Atlas of Language Structures

## Morphology and Syntax

- Word orders freedom and morphology are usually related
- The more freedom in word orders
  - → the less information is conveyed by word positions
  - → the more information should be included in the "symbols"
  - → the richer the morphology
- e.g English vs. Russian (object indicated with -ей):

<div align="center">

*cats eat mice*

Кошки едят мышей

Мышей едят кошки.

Едят кошки мышей.

Едят мышей кошки.

</div>



Constrained word order
Limited or no morphological marking

(Relatively) free word order
Rich morphology

# Semantic Diversity

- Words partition the semantic space
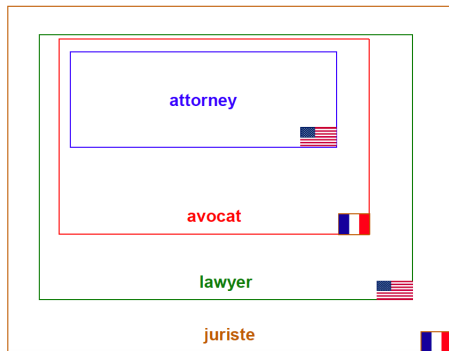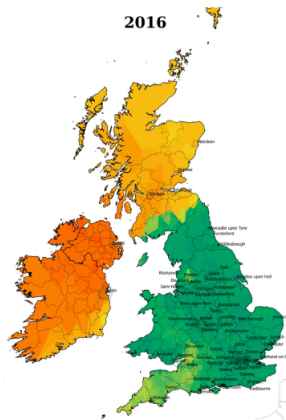- This partition is very diverse across language



Figure: Semantic partitioning between English(US) and French: *laywer* vs *avocat*. (2)

# Variation

- Variation at all level of analysis (phonological, morphological, syntactic, semantic)
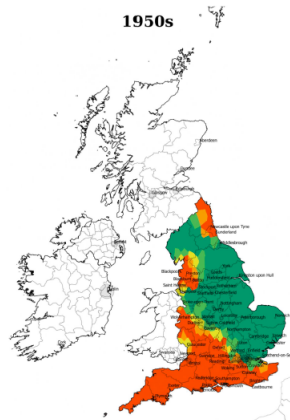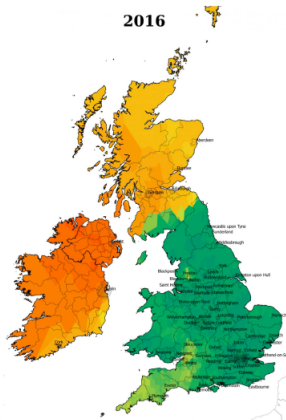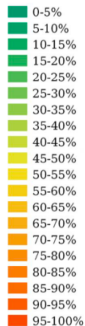- Building NLP with such variance is a great challenge

# Phonetic Variation



Do you pronounce the "r" in "arm" ?

2016

# Phonetic Variation



Do you pronounce the "r" in "arm" ?

- 0-5%
- 5-10%
- 10-15%
- 15-20%
- 20-25%
- 25-30%
- 30-35%
- 35-40%
- 40-45%
- 45-50%
- 50-55%
- 55-60%
- 60-65%
- 65-70%
- 70-75%
- 75-80%
- 80-85%
- 85-90%
- 90-95%
- 95-100%

2016

1950s

# Spelling Variation

anagement   maagement   maanagement
maangement   magagement   magement
mamagement   mamangement   manaagement   manaement
managaement   manageement   manageemnt   managegment
managemaent   managemant   managememt   managemen   managemenet
managementt   managemet   managemetn   managemnent   managemnet
managemnt   managemrnt   managemt   managenent   managenment   managent
managerment   managhement   managmeent   managmgement   managmgment   managnment
manament   manamgement   mananement   manangment   manasgement
manegement   manegment   mangaement   mangagement   mangagment
mangament   mangement   manggement   mangment
mangmt   menagement   mgmt   mgnt
mnagement   mngmnt   mngmt

Figure: Spelling variation of "management" found in Social Media data (2)

# Sociolinguistic Variation



Figure: Non-Canonical Tweet Translated by Bing (2)

# Ambiguity

- Most linguistic observations (speech, text) are open to several interpretation
- We(Humans) disambiguate/find the correct interpretation using all kind of signals (linguistic and extra linguistic)
- Ambiguity can appear at all levels of analysis

# Syntactic Ambiguity



Figure: Syntactic Ambiguity (2)

# Semantic Ambiguity



Figure: Semantic Ambiguity (2)

- Name entity
- Polysemy (man)
- Object/Color (cherry)
- Object/Informal (e.g. the book)

# Ambiguity examples

- Ambiguity! Some examples of ambiguous headlines:

  *Iraqi head seeks arms*

  *Enraged cow injures farmer with axe*

  *San Jose cops kill man with knife*

  *Miners refuse to work after death*

  *Two Soviet ships collide, one dies*

  *Dealers will hear car talk at noon*

- Ambiguity can be lexical, syntactic, pragmatic

# Ambiguity examples

**Human:** Are there direct flights from Paris to Santiago?

        **Bot:** Yes, there is an Air France flight leaving at 11:40PM.

**Human:** How long does it takes to go there?

                **Bot:** The flight takes 14h35m.

**Human:** How much would that cost?

- Needs discourse knowledge, domain knowledge, linguistic knowledge

# Sparsity

Data Sparsity is when many entities (words, morphemes, n-grams, ...) in a corpus have very low observed frequency

Sparsity is the consequence of :

- **Combinatorial** structure of language
  *Combining meaningless sounds into meaningful morphemes or words and meaningful phrases* into sentences. [3]
- **Zipfian** structure of language

NB : Sparsity is one of the greatest challenge of NLP

---
[3] The Origin of Speech, Hockett et. al 1960

# Zipf's law

*Zipf ' s law* can describe many phenomenons of language.

Definition:
$f_w$ frequency of entity w
k frequency rank of entity w

$$f_w(k) \; \alpha \; \frac{1}{k^\theta}$$

Comments

- Zipf law is a Power relation between the rank and frequency
  *The most frequent entities are <u>much</u> more frequent than the less frequent ones*

- Under Zipf law $log(f_w)$ and $log(k)$ are linearly related

# Zipfian structures in Language

Zipf law can be found in many phenomenons in nature.

In Language

- Word frequency
- Syntactic structures frequency

# Zipfian structure of Language



word frequency and rank in *Romeo and Juliet* (linear–linear)

# Zipfian structure of Language : Lexicon



$$f(rank; s, N) = \frac{1/rank^s}{\sum_{n=1}^{N} 1/n^s}$$

Figure: rank vs frequency for the top 10M words in wikipedia

# Zipfian structure of Language : Syntax



Figure: rank vs frequency for automatically parsed corpus

# Zipf's law and Sparsity

- The Zipf's law is a long tail distribution
- Many entities (words, syntactic structure,...) have a very low frequency

  $\rightarrow$ sparsity

# What is Natural Language Processing ?

# What is Natural Language Processing ?
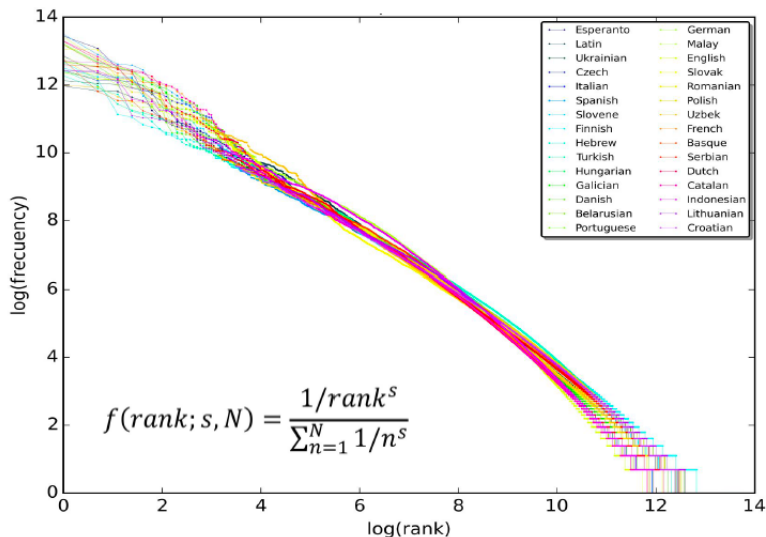
- Process, analyze and/or produce natural language
- Interact with computers using natural language

- Natural language 'understanding':
  - language as input $\mapsto$ "information" as output
- Natural language generation:
  - "information" as input $\mapsto$ language as output
- Sometimes, both: machine translation, summarization, question answering

- Strongly related fields:
  - machine learning,
  - artificial intelligence,
  - deep learning
  - (computational) linguistics

# What is Natural Language Processing ?

In a nutshell, NLP consists in handling the complexities of language systematically "to do something"

- Raw Text $\rightarrow$ Structured Information
- Raw Text $\rightarrow$ Controlled Text

# Brief History of NLP



| 1970 | 2000 | 2012 |

**Symbolic approaches**

•**Computational expertise:**
Formal grammars (algebraic grammars, mildly context-sensitive grammars, polynomial languages…), parsing algorithms, dynamic programming

•**Comp. linguistics expertise:**
Formal and descriptive linguistics, grammar engineering, development of lexical resources

**Statistical approaches**

•**Computational expertise:**
(statistical) machine learning, supervised, semi-supervised and non-supervised (PCFG, CRF, MEMM, discriminative algorithms…), hybrid approaches

•**Comp. linguistics expertise:**
development of annotated corpora (training dataset), development of lexical resources

**Neural approaches**

•**Comp. expertise:**
neural networks, deep learning, end-to-end training

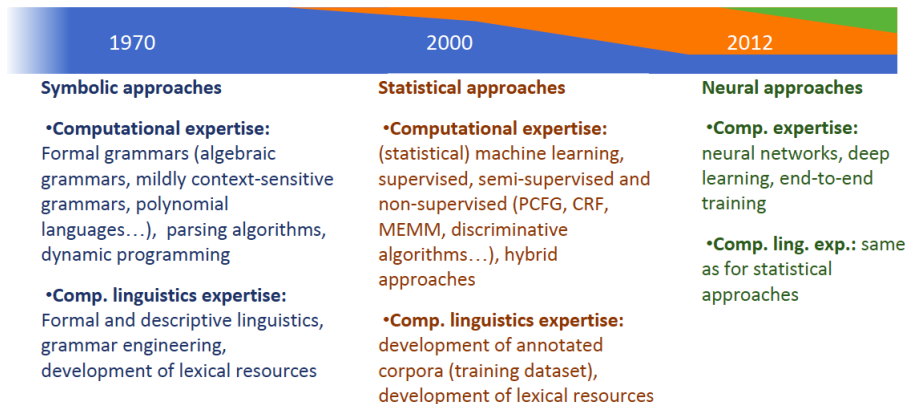•**Comp. ling. exp.:** same as for statistical approaches

Figure: Brief history of NLP (2)

# NLP prediction pipeline

# NLP training pipeline

# NLP in the real-world



- Building NLP systems is an iterative cycle...
- Composed of **Human & Machine Learning**

# Outline of the course

1. The Why and What of Natural Language Processing
2. Representing text with vectors
3. Modeling textual data
4. Neural Natural Language Processing
5. Language Modelling
6. Transfer Learning with Neural Modeling for NLP

# References I

[Kracht] Kracht, M. Introduction to linguistics.

[2] Sagot, B. (2019). Algorithms for speech and natural language processing, lectures ens-saclay.