# Natural language processing
## Introduction and course overview

X. Hilaire

ESIEE Paris, IT department
x.hilaire@esiee.fr
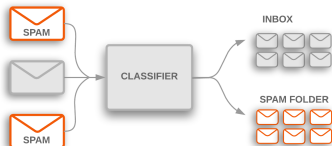
August 27, 2023

# Outline

# What natural language processing (NLP) ?

- In a nutshell: NLP is a field of computer science which aims to reproduce everything that humans know to do with language by a computer.
- Includes (but not limited to) various tasks such as :
  - Text classification [1]
  - Text clustering and topic detection [1]
  - Sentiment and opinion analysis [1]
  - Text summarization
  - Machine translation [1]
  - Text generation
  - Question answering (ChatGPT)

---

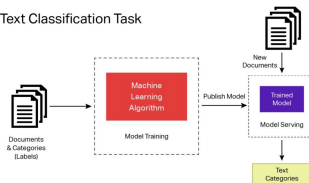[1]Problem tackled in the course, at least once in a lab

# Text classification

**Text classification** is to label a given text using predefined labels. Examples: news, finance, sports, cars, internet, etc. A simple, yet useful example is SPAM/NON SPAM classification.



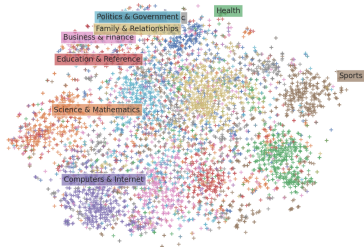---

[2]https://towardsdatascience.com
[3]https://imerit.net

# Text clustering

- In essence, text clustering is the unsupervised version of text classification
- Goal = to segment a text or corpus in different parts, but those don't have any predefined label



4

A reference library for text clustering: *Carrot*[2]
https://github.com/carrot2/carrot2
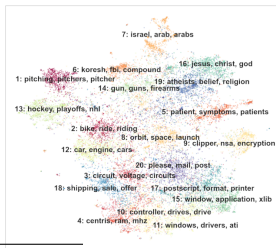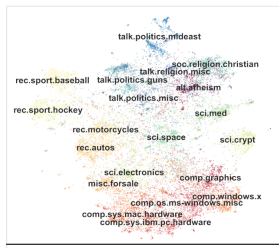https://search.carrot2.org/

---

[4] ZSL Blog Post

# Topic detection and modeling

- Closely related to text clustering
- A **topic** is a set of words which are highly correlated one to each other (a "strong" cluster)
- Assuming a document is related to several topics:

$$p(word_i|doc_j) = \sum_k p(word_i|topic_k)p(topic_k|doc_j)$$

This bayesian expansion turns out to work better than working on words directly.
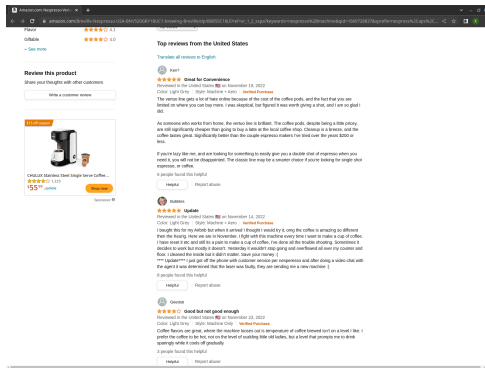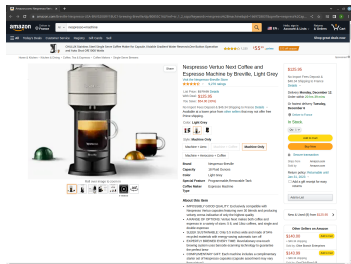


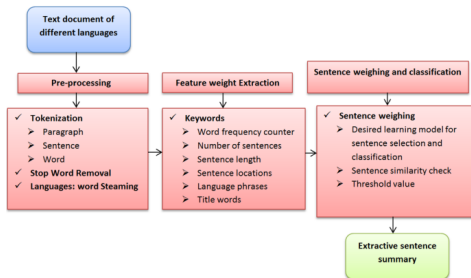[5]sbert.net

# Sentiment and opinion analysis

Extracting subjective information from text :

- Has a given text a positive / negative connotation ?
- Is it enthusiastic / critical, joyful / sad, interesting / not interesting ?
- Has a given product received a good / bad review ?

# Text summarization

- Aim = to compress a text while minimizing the loss of its meaning
- Very challenging task, as it requires a deep understanding of the text
- Two broad families of methods : extractive and abstractive
- Extractive methods
  - Extract important sentences from the text, and drop others.
  - A scoring function is needed to rank the importance of sentences. LSA (Latent semantic analysis), which we will see, is often used to do so.



*"Study of automatic text summarization approaches in different languages"* [2]

# Text summarization

- Abstractive methods
  - Analyze the semantic information of every sentence of the text
  - Produce shorter, completely new sentences having the same meaning
- Two representative methods of this family are that of:
  - Nallapati et al. (IBM) "*Abstractive Text Summarization Using Sequence-to-Sequence RNNs and Beyond*" [3]
  - See et al. (Google) "*Get To The Point: Summarization with Pointer-Generator Networks*" [5] : both use sequence-to-sequence alignments of words.
- For a survey, see Li Jinpeng et al. "*Survey on Automatic Text Summarization*" [1]

# Text summarization



Feature-rich-encoder.
From Nallapati et al. [3]



Pointer-generator model.
From See et al. [5]

## Machine translation

- Unmissable task, as very useful for centuries.
- Difficult task, as even human translators can disagree (*traduttore, traditore*)
- Very old problem, plethora of methods and products available (DeepL, Systran, Google translate, Microsoft Translator, Reverso, ...)
- For a recent survey, see, e.g. I. Rivera-Trigueros "*Machine translation systems and quality assessment: a systematic review* " [4]

From `https://nlp.johnsnowlabs.com`

## The most widely used NLP library in the enterprise

Which NLP libraries does your organization use [check all that apply]?



Spark NLP — 33%
spaCy — 26%
Allen NLP — 23%
nltk — 21%
Stanford CoreNLP — 20%
Gensim — 18%
Hugging Face — 18%
Other — 11%
Rasa NLU — 10%

Source: 2020 NLP Industry Survey, by **Gradient Flow**.

# Outline

## Objectives

To know the most important algorithms used in NLP :

1. Text preprocessing (tokenization) ☞ painful but mandatory step, needed to :
   - extract the textual parts of documents, often much richer (LaTeX, HTML, XML, etc.)
   - normalize words, possibly correcting typos on-the-fly

2. Word embeddings ☞ maps words from a sparse space to a dense space. Actually <u>very</u> important :
   - moves words which are semantically close to close representations in the dense space
   - allows to substitute a word by a synonym with minimal distortion

3. Topic detection, text classification, sentiment analysis ☞ old technique, borrowed from dimensionality reduction.

4. Machine translation ☞ automatically translates a text from one language to another with minimal semantic distortion

5. Transformers are not covered in this course due to lack of time, but in AIC-5201C "Deep-Learning and its Applications"

# Course overview

The course heavily relies on some fundamental tools :

- Singular value decomposition (SVD) for dimensionality reduction and topic detection
- Fisher kernel linear discriminant analysis (k-LDA) and SVM for text classification and sentiment analysis
- Ordinary (non-recurrent) neural network and **backpropagation** $\leftarrow$ very important
- Elman networks, recurrent neural networks (RNN)
- Long short term memory (LSTM) and gated recurrent units (GRU) networks

Don't underestimate them !

## Course overview

Schedule:

- Lectures and exercise sessions / labs are alternated
  - about 20–40mn of lecture
  - then one exercise on the topic seen
- Do not pay strict attention to your schedule on https://planif.esiee.fr : the C/TD/TP breakdown is only meaningful in terms of accounting

About labs:

- Programming language is Python (v3)
- The NLP part involves NLTK, Gensim, and basic Python libraries (numpy, scikit-learn, matplotlib, ...)
- Neural networks can be implemented with Keras (simple, but slow) or PyTorch (more complex, but faster). TensorFlow allowed (for efficiency reasons), but beware of complex API, and expect limited support from me.

# Overview

- Labs have been validated on ESIEE's machines, Linux OS (Debian 11.5)
- You can work on your own laptop or environment, but the final lab you will submit <u>must</u> work on ESIEE's machines, where it will be assessed
- RAM consumption is limited, but training is CPU intensive

# Grading policy

- Final grade = 10% lab2 (sentiment analysis) + 40% lab4 (machine translation) + 50% written exam
- Labs must be done by pairs. Instructions for declaring pairs and submitting will be provided by mail, and involve my server `https://mvproxy.esiee.fr`
- The written exam has two parts:
    - First is 20–30mn, close book (no documents), rewarded 4–5 points. It typically consists of short questions, testing your knowledge about the lectures.
    - Second is 1h30–1h40mn, open book (all documents are allowed), rewarded 15–16 points. It typically consists of case studies.
- Electronic devices of all kinds, are forbidden during the written exam.
- Should you be granted a reset / catch up exam, the conditions remain exactly the same.

## Documents, post-assistance

Course web page:

- The homepage of the course is
  https://perso.esiee.fr/~hilairex/AIC-5102B
- All the information and documents you need will be added to it on
  the fly

Post-assistance:

- I am usually available on Thursday afternoons (1pm-7pm) at my
  office 5352 for post-assistance
- If you have an urgent problem, you can always reach me by email:
  x.hilaire@esiee.fr

📄 Li Jinpeng, Zhang Chuang, Xiaojun Chen, Yue Hu, and Pengcheng Liao.
Survey on automatic text summarization.
Journal of Computer Research and Development, 58(1):1, 2021.

📄 Yogesh Kumar, Komalpreet Kaur, and Sukhpreet Kaur.
Study of automatic text summarization approaches in different languages.
Artif. Intell. Rev., 54(8):5897–5929, 2021.

📄 Ramesh Nallapati, Bowen Zhou, Cicero Nogueira dos santos, Caglar Gulcehre, and Bing Xiang.
Abstractive text summarization using sequence-to-sequence rnns and beyond.
2016.

📄 Irene Rivera Trigueros.
Machine translation systems and quality assessment: a systematic review, 4 2021.

📄 Abigail See, Peter J. Liu, and Christopher D. Manning.

Get to the point: Summarization with pointer-generator networks, 2017.