

# 91258 - Natural Language Processing

## Lesson 1. Introduction

Alberto Barrón-Cedeño

Alma Mater Studiorum-Università di Bologna  
a.barron@unibo.it @albarron\_

02/10/2023



## Table of Contents


Materials

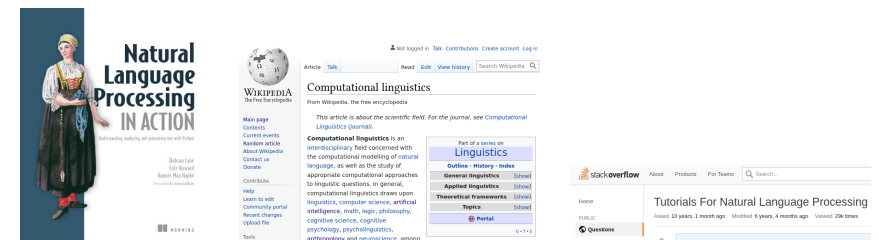
Introduction

Requirements

## Materials






## Core Bibliography

1. Lane et al. (2019)'s  **Natural Language Processing in Action**<sup>1</sup>
2. Numerous **Wikipedia articles** on relevant topics
3. Multiple online forums



<sup>1</sup><https://www.manning.com/books/natural-language-processing-in-action>

## Complementary Bibliography

1. Intro to computing for text  
 K.W. Church's **Unix for poets**<sup>2</sup>
2. For social media analysis  
 Hovy (2021)'s **Text Analysis in Python for Social Scientists**<sup>3</sup>
3. A basic intro in Italian  
 Nissim and Pannitto (2022)'s **Che cos'è la linguistica computazionale**
4. From linguistics  
 Bender (2013)'s **Linguistic fundamentals for natural language processing: 100 essentials from morphology and syntax**<sup>4</sup>
5. Advanced  
 Goldberg (2017)'s **Neural Network Methods for NLP**<sup>5</sup>

<sup>2</sup><https://web.stanford.edu/class/cs124/kwc-unix-for-poets.pdf>


<sup>3</sup><https://doi.org/10.1017/9781108873352>

<sup>4</sup><https://doi.org/10.2200/S00493ED1V01Y201303HLT020>

<sup>5</sup><https://doi.org/10.2200/S00762ED1V01Y201703HLT037>

## Lesson coordinates

Slides, code, and more are all available at:

 [albarron.github.io/teaching/natural-language-processing](https://albarron.github.io/teaching/natural-language-processing)



## Tools

### Essential

Python 3 development framework on any modern OS

1. Command line **or**
2. Integrated development Environment; e.g., Pycharm<sup>6</sup>, Eclipse<sup>7</sup> **or**
3. Jupyter notebook; e.g., Google's colab<sup>8</sup>, local Jupyter<sup>9</sup>

### Desirable<sup>10</sup>

1. Git Version control system; e.g.,  Gitlab<sup>11</sup> **or**  Github<sup>12</sup>
2.  $\text{\LaTeX}$  system for document preparation

<sup>6</sup><https://www.jetbrains.com/pycharm/>

<sup>7</sup><https://www.eclipse.org/>

<sup>8</sup><https://colab.research.google.com/>

<sup>9</sup><https://jupyter.org/>

<sup>10</sup>Could be part of "Advanced Research Skills Lab"

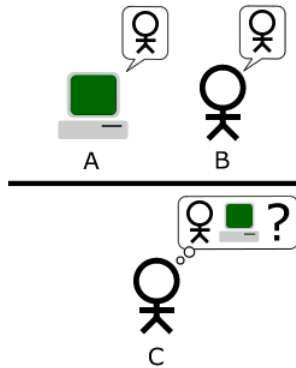
<sup>11</sup><https://gitlab.com>


<sup>12</sup><https://github.com>

## Introduction

## Introduction

Natural language as a measure of intelligence



 Turing (1950). "Computing machinery and intelligence". Mind. 59(236)

[upload.wikimedia.org/wikipedia/commons/e/e4/Turing\\_Test\\_version\\_3.png](https://upload.wikimedia.org/wikipedia/commons/e/e4/Turing_Test_version_3.png)

## Introduction

CL vs NLP

### Computational linguistics<sup>13</sup>

- **Interdisciplinary** field concerned with the **computational** (it used to say "statistical or rule-based"! ) **modeling of natural language**
- Study of appropriate computational approaches to **linguistic questions**

### Natural Language Processing<sup>14</sup>

- Subfield of **computer science** and **linguistics** [...] concerned with giving computers the ability to support and manipulate speech
- Processing natural language datasets, such as text corpora or speech corpora, using either rule-based or probabilistic machine learning approaches

<sup>13</sup>[https://en.wikipedia.org/wiki/Computational\\_linguistics](https://en.wikipedia.org/wiki/Computational_linguistics)

<sup>14</sup>[https://en.wikipedia.org/wiki/Natural\\_language\\_processing](https://en.wikipedia.org/wiki/Natural_language_processing)

## Introduction

CL vs NLP

### Natural Language Processing (Lane et al., 2019, p. 4)

- Area of research in computer science and artificial intelligence concerned with **processing natural languages**
- This processing generally involves **translating natural language into data** (numbers) that a computer can use to learn about the world

The term **natural language processing** is nowadays considered to be a near-synonym of **computational linguistics** and (human) **language technology**.<sup>15</sup>

<sup>15</sup>[https://en.wikipedia.org/wiki/Computational\\_linguistics](https://en.wikipedia.org/wiki/Computational_linguistics)

## Introduction

### Rule-based vs Statistical NLP

## Introduction

### Rule-based NLP

Models are based on a number of hand-crafted rules or grammars



Diagram borrowed from L. Moroney's Introduction to TensorFlow for Artificial Intelligence, Machine Learning, and Deep Learning

## Introduction

### Rule-based NLP

Models are based on a number of hand-crafted rules or grammars

```
greeting_inputs = ("hey", "morning", "evening", "hi",  
                  "whatsup", "hello")  
greeting_responses = ["hey", "hey hows you?", "*nods*",  
                      "hello, how you doing", "hello",  
                      "Welcome, I am good and you"]
```

```
def generate_greeting_response(input):  
    for token in input.split():  
        if token.lower() in greeting_inputs:  
            return random.choice(greeting_responses)
```

Derived from <https://stackabuse.com/python-for-nlp-creating-a-rule-based-chatbot/>

[//stackabuse.com/python-for-nlp-creating-a-rule-based-chatbot/](https://stackabuse.com/python-for-nlp-creating-a-rule-based-chatbot/)

## Introduction

### Statistical NLP

Models are tuned on *annotated* data

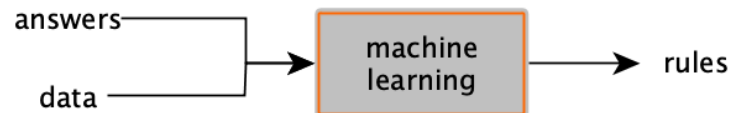
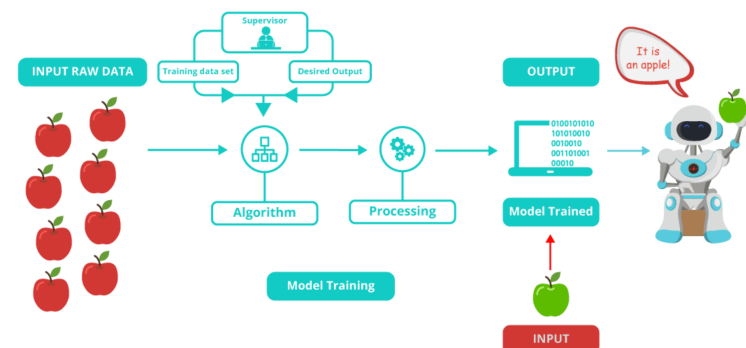


Diagram borrowed from L. Moroney's Introduction to TensorFlow for Artificial Intelligence, Machine Learning, and Deep Learning

## Introduction

### Statistical NLP

Models are tuned on annotated data

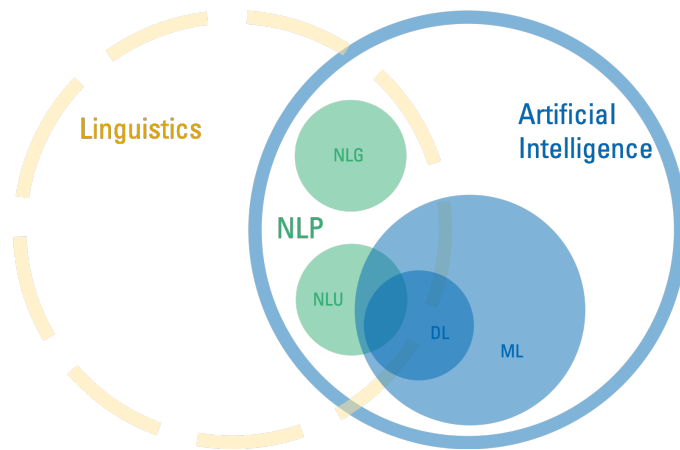


Borrowed from

<https://www.edureka.co/blog/machine-learning-tutorial>

## Introduction

The NLP neighborhood



Borrowed from <https://www.retresco.de/en/how-to-ai-natural-language-processing/>

## Requirements & Evaluation

## Introduction

Non-exhaustive list of NLP applications with examples

🔍 <b>Search</b>	web search engines · text autocompletion
✍️ <b>Editing</b>	grammar issues identification
💬 <b>Dialogue</b>	chatbot creation
✉️ <b>Email</b>	spam filtering · message classification
📄 <b>Text mining</b>	(multi-)document summarisation
📰 <b>News analysis</b>	event identification · fact checking
👤 <b>Forensics</b>	plagiarism detection · authorship attribution
👍 <b>Sentiment analysis</b>	product review ranking · opinion mining
✍️ <b>Creative writing</b>	text generation with a narrative and style
🗣️ <b>Translation</b>	translation · quality estimation

Partially derived from (Lane et al., 2019, p. 8)

## Requirements

### Necessary

- ▶ Linguistics
- ▶ Algebra
- ▶ Programming in **Python**

### Desirable

- ▶ Intermediate programming (e.g., object-oriented, testing)
- ▶ High-performance computing (e.g., slurm)<sup>16</sup>

<sup>16</sup>Could be part of “Advanced Research Skills Lab”

## Evaluation: **One final project**

You will address a relevant problem...

- ▶ within the range of your own (research) interests
- ▶ participating (formally) in a shared task
- ▶ proposed by me, if you prefer

## Evaluation: **One final project**

### Typical pipeline

1. You propose a topic/problem. We assess if it is reasonable, doable. . .
2. You compile data, study the problem, design experiments, code. . .

### **IF you plan for a publication**<sup>17</sup>

- ▶ We meet regularly to see the advances and shape the experiments, submissions, and/or paper towards the submission deadline

### **ELSE**

- ▶ We could meet sporadically, if you need it

3. You submit a written report (~ 7 pages) **1 week before the appello**
4. We meet on the date of the appello to discuss about your project, in the context of the lecture

---

<sup>17</sup>Talk to me well in advance; it would require my heavy involvement to target a high quality

## Evaluation: **Final mark**

Combination of the quality of the experiments, report, code, and oral discussion

### **Targetting 30L?**

If I let you submit a paper, it is very likely. But it is not the only way. . .

$$p(30L \mid \text{paper submitted} == \text{True}) \approx 0.85 \quad (1)$$

$$p(30L \mid \text{paper submitted} == \text{False}) \approx 0.15 \quad (2)$$

## Evaluation: Previous final projects

---

### **2022–2023**

- 🗣️ Sentiment analysis of video game reviews
- ⚖️ Authorship attribution: machine vs human

### **2021–2022**

- ✖️ Hate Speech Detection in Incel Online Spaces
- ♂️ Fishing for catfishes: predicting the author gender in Reddit

---

\* student with previous programming skills

- turned into (part of a) thesis    ♠️ turned into a publication

## Evaluation: Previous final projects

### 2020–2021

- ✂ Semantic similarity between originals and machine translations •
- 🔑 Definition extraction on food-related Wikipedia articles •
- ☰ Identifying Characters' Lines in Original and Translated Plays •
- 🐦 Classifying an Imbalanced Dataset with CNN, RNN, and LSTM

### 2019–2020

- ❤ AriEmozione: Identifying Emotions in Opera Verses \*♣
- 🐦 UniBO@AMI: A Multi-Class Approach to Misogyny and Aggressiveness Identification on Twitter Posts Using ALBERTo \*♣

\* students with previous programming skills

- turned into (part of a) thesis ♣ turned into a publication

Visit the **projects section** of the class website for details, reports and papers

## References

Bender, E. M.

2013. *Linguistic Fundamentals for Natural Language Processing: 100 Essentials from Morphology and Syntax*. Morgan & Claypool Publishers.

Goldberg, Y.

2017. *Neural Network Methods in Natural Language Processing*. Morgan & Claypool Publishers.

Hovy, D.

2021. *Text Analysis in Python for Social Scientists: Discovery and Exploration*, Elements in Quantitative and Computational Methods for the Social Sciences. Cambridge University Press.

Lane, H., C. Howard, and H. Hapkem

2019. *Natural Language Processing in Action*. Shelter Island, NY: Manning Publication Co.

Nissim, M. and L. Pannitto

2022. *Che cos'è la linguistica computazionale*. Carocci editore.