

Natural Language Processing (NLP)

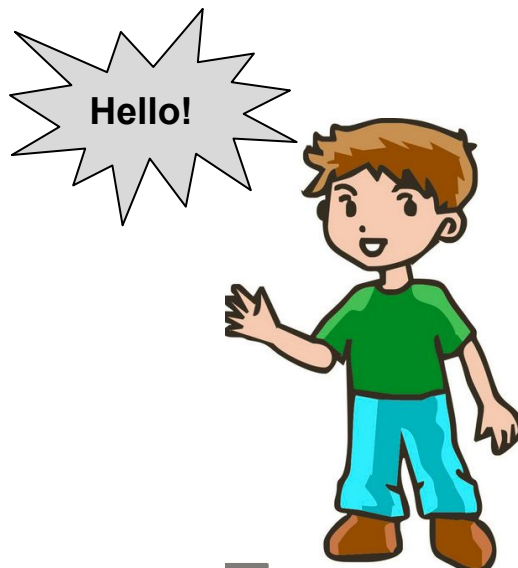
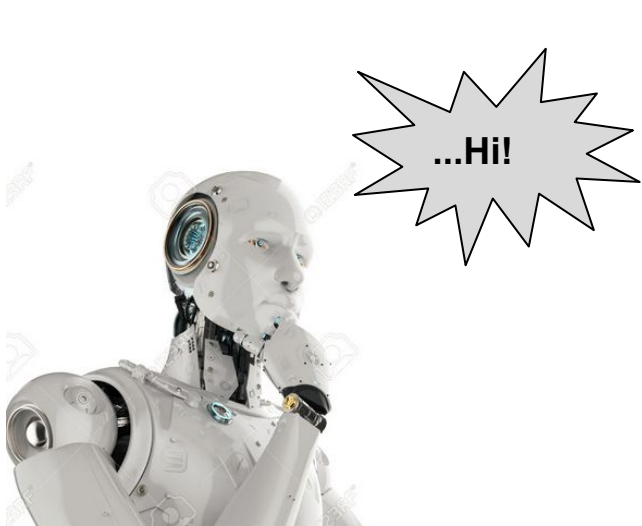
by S.B., M.Z.

Outline

- ❖ NLP Intro
- ❖ Examples
- ❖ Basic Commands
 - Tokenization
 - Stemming & Lemmatization
 - Stop Words
 - POS-Tag
- ❖ Advanced
 - Sentiment Analysis
- ❖ Limitations & Future
- ❖ Conclusion

What is NLP?

Definition: “a subfield of linguistics, computer science, information engineering, and artificial intelligence where the goal is to be able to get machines (computers) to understand the meaning of human (natural) language.” ([Wikipedia](#))



High Level Structure

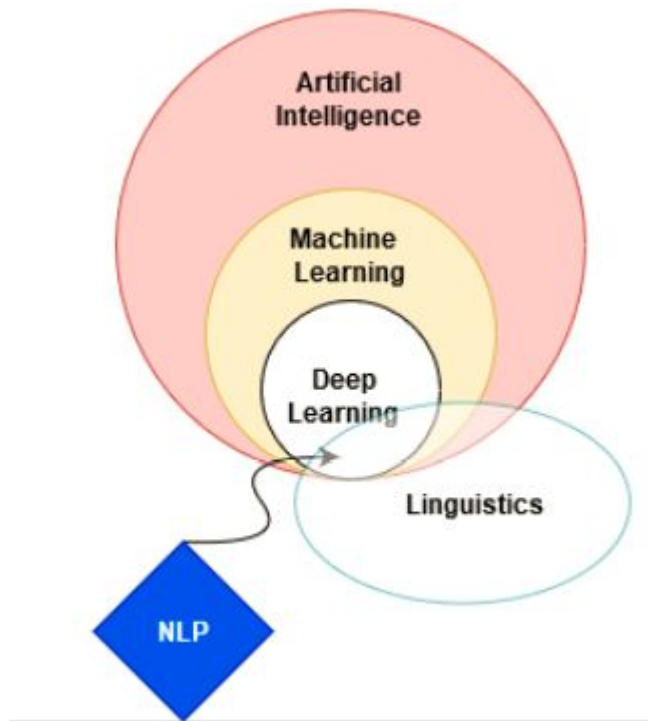
Artificial Intelligence: concerned with making computers intelligent (less reliant on human instruction).

Machine Learning: set of tools for making inference and predictions (computer science and statistics).

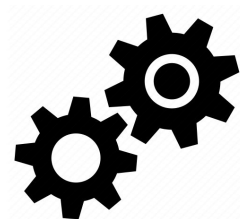
Deep Learning: special area of ML where a lot more data is required to establish relationships strengths (text/pics) uses multiple layers to progressively extract higher level features from the raw input (edges» shape » identify object).

- ❖ We have complex problems (language)
- ❖ Lots of data that we can feed into the model

Linguistics: study of language through sound and meaning.



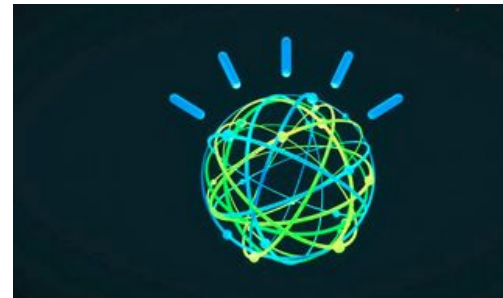
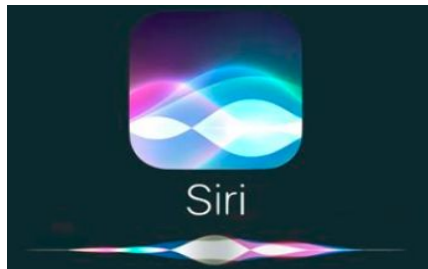
Why is it important?



- ❖ NLP processing helps computers communicate with humans in their own language and scales other language-related tasks.
 - read text, hear speech, interpret it, measure sentiment and determine which parts are important.
- ❖ Machines can analyze more language-based data than humans, without fatigue and in a consistent, unbiased way.
 - More efficient & less biased

Problems Solved with NLP

❖ Voice Driven Interfaces



Problems Solved with NLP



❖ Machine Translation

ITALIAN - DETECTED ITALIAN ENGLISH GREEK ↕ ENGLISH ITALIAN GREEK

Anche io, nei piu' soclusi posti della mia mente, sono capace di follia. |

I, too, in the most sociable places of my mind, am capable of madness.

73/5000

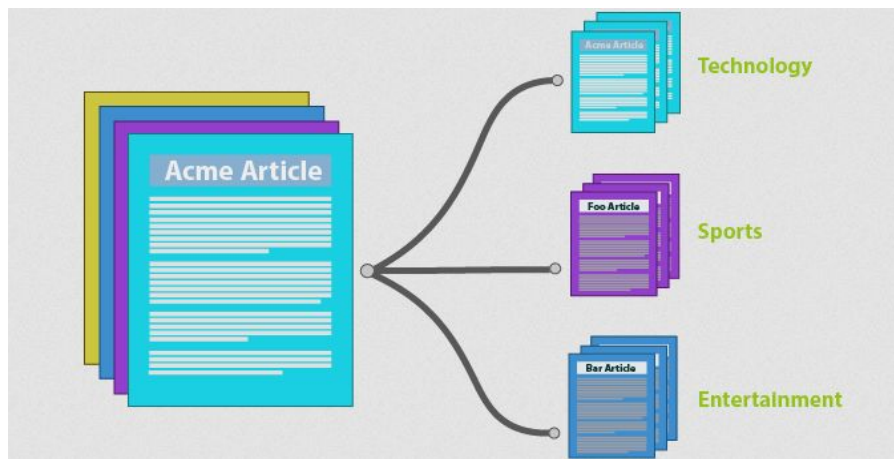
Anche io, nei piu' **socclusi** posti della mia mente, sono capace di follia.

SECLUSI

I, too, in the most **soccluded** places of my mind, am capable of madness.

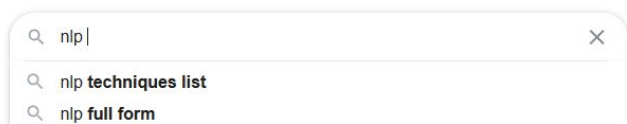
Problems Solved with NLP

❖ Document Categorization

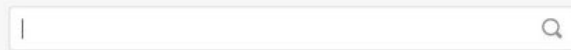


Problems Solved with NLP

❖ Search Engines



DuckDuckGo



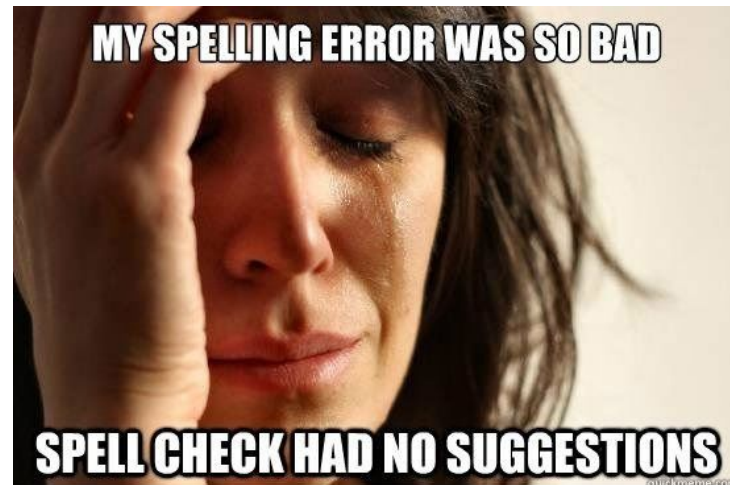
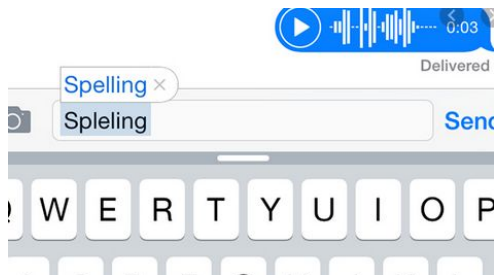
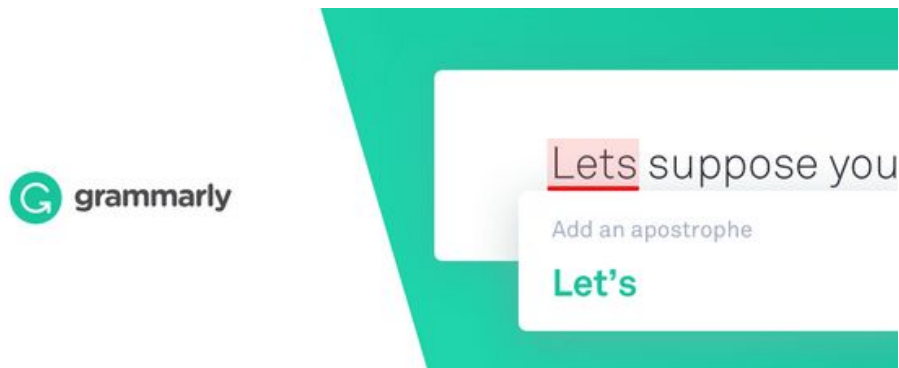
The search engine that doesn't track you. [Help Spread DuckDuckGo!](#)

many more....



Problems Solved with NLP

- ❖ Spell Checker (*autocomplete, autocorrect*)

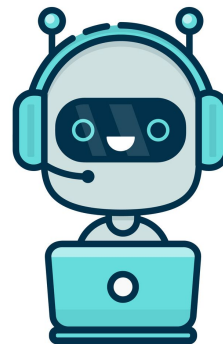
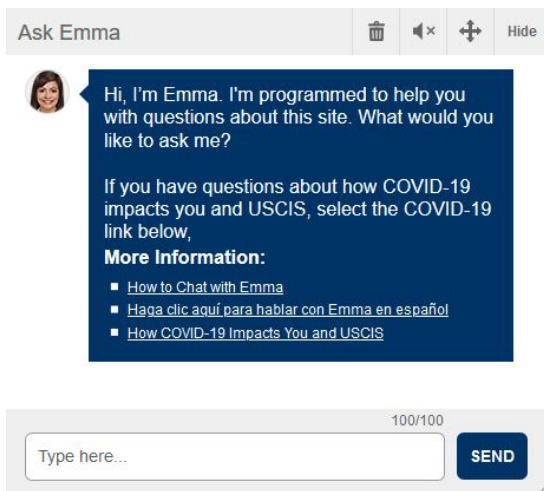


more funny stuff at FUNNYASDUCK.NET



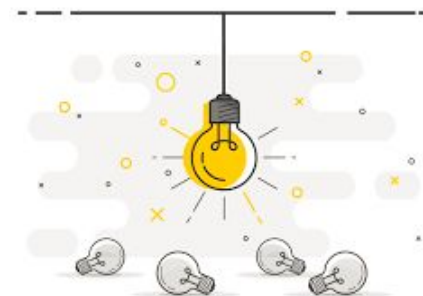
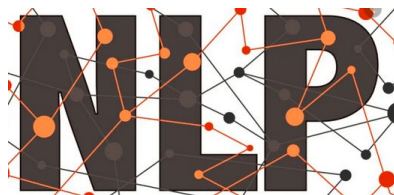
Problems Solved with NLP

❖ Robot Chat



Major Pillars of NLP

- ❖ **Content categorization.** A linguistic-based document summary, including search and indexing, content alerts and duplication detection.
- ❖ **Topic discovery and modeling.** Accurately capture the meaning and themes in text collections, and apply advanced analytics to text, like optimization and forecasting.
- ❖ **Contextual extraction.** Automatically pull structured information from text-based sources.
- ❖ **Sentiment analysis.** Identifying the mood or subjective opinions within large amounts of text, including average sentiment and opinion mining.
- ❖ **Speech-to-text and text-to-speech conversion.** Transforming voice commands into written text, and vice versa.
- ❖ **Document summarization.** Automatically generating synopses of large bodies of text.
- ❖ **Machine translation.** Translating file/voice input from one language to the other



Basic Functions

- ❖ Packages
- ❖ Tokenization
- ❖ Stemming
- ❖ Lemmatization
- ❖ Stop Words
- ❖ Sentiment Analysis

Packages

- Data Frame Manipulation (selection, deletion, sorting...)
 - Numpy
 - Pandas
- Data Processing (specific operation pertinent to NLP procedures)
 - NLTK
 - TextBlob
- Modeling
 - Sklearn
 - KNN
 - Linear Regression

Tokenization

“The process of segmenting running text into sentences and words. It is the task of cutting a text into pieces called **tokens**.” Such tokens can vary according to user needs: split text into individual words, sentences, or **n-grams** (collections of n words).

Example:

Word tokenization: “Hello Everyone” ➞ “Hello”, “Everyone”

Sentence tokenization: “ Hello Everyone. How are you all doing? I hope you are well.”
➞ “Hello Everyone. ”, “How are you all doing?”, “ I hope you are well.”

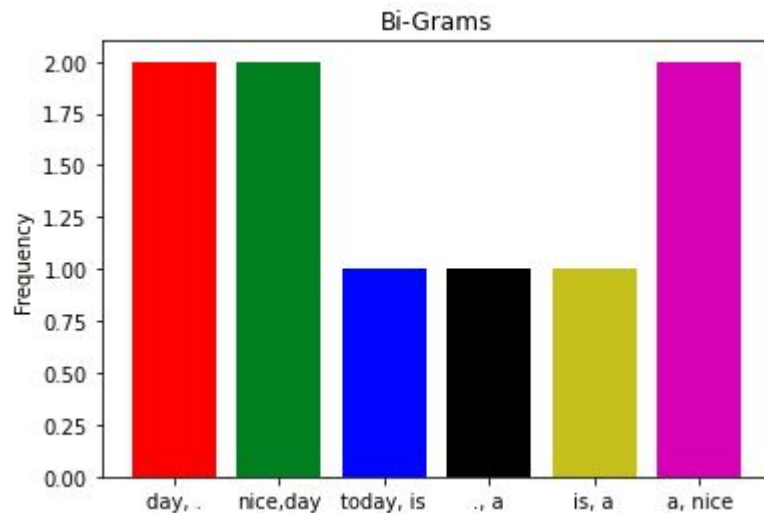
N-grams

“N-grams can be considered as a special type of token ”

Example: **Bi-Gram** : “Today is a nice day. A nice day. ”

➤ “Today is”, “ is a ”, “a nice”, “nice day”, ‘day .’, ‘. A’, ‘A nice’, ‘nice day’, ‘day .’

Today is	is a	a nice	nice day	day .	. A
1	1	2	2	2	1



Tokenization ...

```
#Text
text = 'Welcome to NextGen seminar. I hope you are enjoying learning about NLP.'
#Function sent_tokenize(arg) which splits the text into sentences
sentence_token = sent_tokenize(text)
#Retruns a list object
print(type(sentence_token))

#Display sentences
for index in range(len(sentence_token)):
    print("Sentence ",index+1, ": ",sentence_token[index])

<class 'list'>
Sentence 1 : Welcome to NextGen seminar.
Sentence 2 : I hope you are enjoying learning about NLP.
```

Stemming

“Stemmers remove morphological **affixes** from words, leaving only the word stem.” In short, removing ‘*extras*’ from a word and reduce it to its basic root.

Example: beautiful ➤ beauty , dies ➤ die, flying ➤ fly, mules ➤ mule

OverStemming: stemming a word when non-necessary (frequent problem).

Universal ➤ Universe , University ➤ Universe

Porter Stemmer: common stemmer used, although it is not the most precise, better stemmers exist, but it is useful for basic stemming needs.

Stemming ...

```
#Create stemmer object
stemmer = PorterStemmer()

#Define a word
single_word= "Programers"

#Stem word by calling the stemmer object and calling the function stem: stemmer.stem()
stem_single_word = stemmer.stem(single_word)

#Print results
print("Word: ",single_word," Stemmmmed version: ", stem_single_word)
```

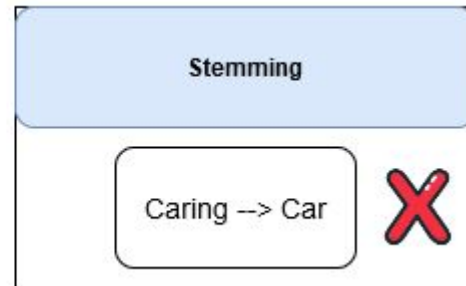
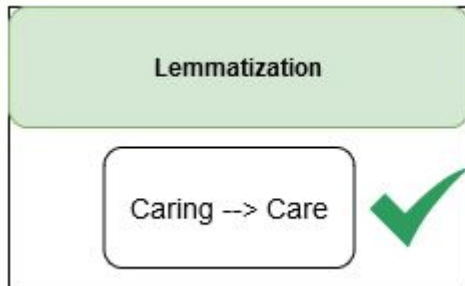
```
Word:  Programers  Stemmmmed version:  program
```

Lemmatization

Similar task to stemming, in the sense that we are modifying a word to its root, but with the difference that in this scenario the change is based on “**meaning**” rather than affixes. Similar words that have similar meaning are considered equal.

Example: verbs in the past tense are changed to the present (went ➤ go), synonyms are unified (best ➤ good).

Efficiency of this libraries is based on the dictionaries (vocabulary) that they use.



Lemmatization...

```
#Lemmatization: applying universal meaning to similar words  went → go, best → good
#Create a Lemma object
lemmatizer = WordNetLemmatizer()

#Define a word
single_word= "went"

#Lemmatize word by calling the lemma object and calling the function lemmatize:  lemmatizer.lemmatize()
lemma_single_word = lemmatizer.lemmatize(single_word, pos='v')

#Print results
print("Word: ",single_word," Lemmatized version: ", lemma_single_word)
```

```
Word:  went  Lemmatized version:  go
```

Stop Words

“Stop Words are words which **do not** contain important significance to be used in search queries. Usually, these words are filtered out from search queries because they return a vast amount of unnecessary information.

Example: words that are commonly used in the English language such as 'as, the, be, are' etc.

Each programming language will give its own list of stop words to use, and you can always add more stop words, or create your own customized version.

Stop Words...

```
#Stop Words: meaningless words that appear in text but that give/or take no additional meaning from major idea
from nltk.corpus import stopwords

#Select stop words for English language
stop = stopwords.words('english')

#Check what the type of the object stop is, its length and contents
print(type(stop))
print(len(stop))
print(stop)
```

```
<class 'list'>
```

```
179
```

```
['i', 'me', 'my', 'myself', 'we', 'our', 'ours', 'ourselves', 'you', "you're", "you've", "you'll", "you'd", 'your', 'yours', 'yourself', 'yourselves']
```

Sentiment Analysis

“Sentiment analysis is a machine learning technique that detects polarity (e.g. a *positive* or *negative* opinion) within text, whether a whole document, paragraph, sentence, or clause.”

Examples

```
good = TextBlob("I love you.")  
bad = TextBlob("I hate you.")  
neutral = TextBlob("Maybe I like you, maybe I don't")  
print(text.polarity)  
print(bad.polarity)  
print(neutral.polarity)
```

0.5

-0.8

0.0

Sentiment Analysis...

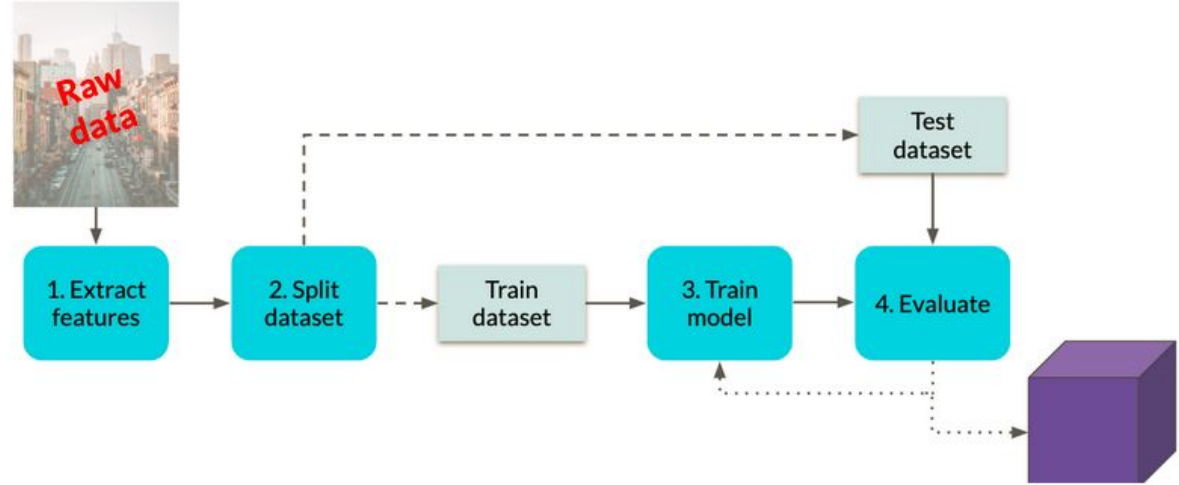
```
#Sentiment Analysis: determining the positivity/negativity of a determinate token (word/sentence)
text = TextBlob("Python is horrible! Hello. You are beautiful. I AM NOT SURE. I am very sure. What a day!")

num=1
for sentence in text.sentences:
    print(num, " ", sentence)
    if(sentence.sentiment.polarity>0):
        print('    Score: positive')
    elif(sentence.sentiment.polarity==0):
        print('    Score: neutral')
    elif(sentence.sentiment.polarity<0):
        print('    Score: negative')
    print("\n")
    num +=1
```

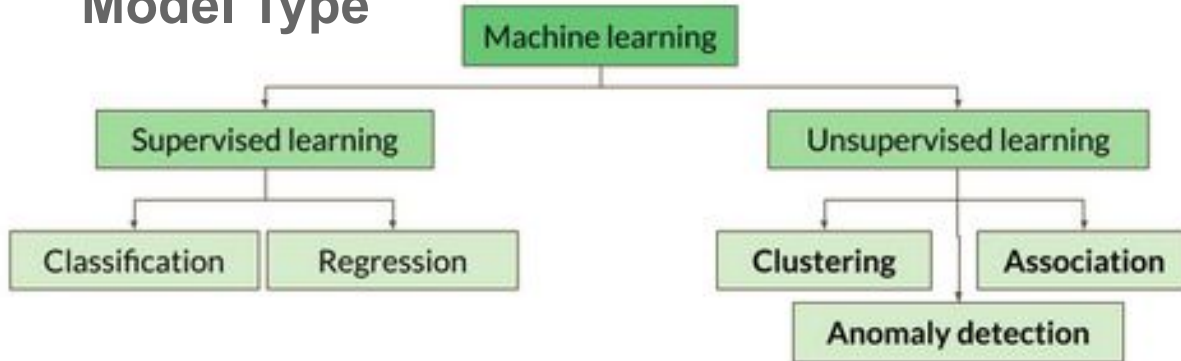
1	Python is horrible! Score: negative	2	Hello. Score: neutral	3	You are beautiful. Score: positive
4	I AM NOT SURE. Score: negative	5	I am very sure. Score: positive	6	What a day! Score: neutral

ML & NLP

Machine learning workflow



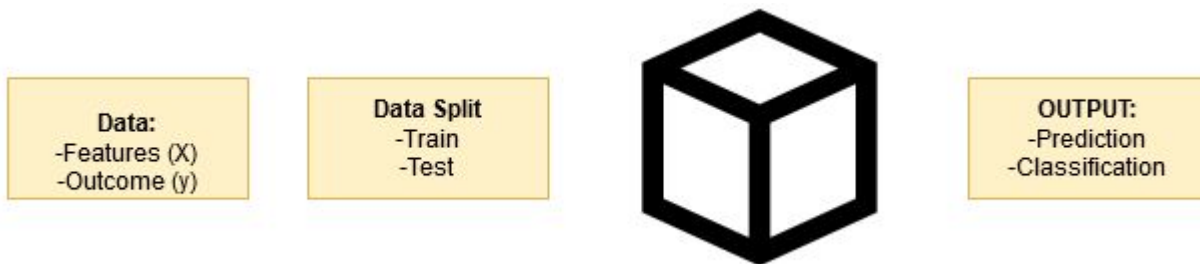
Model Type



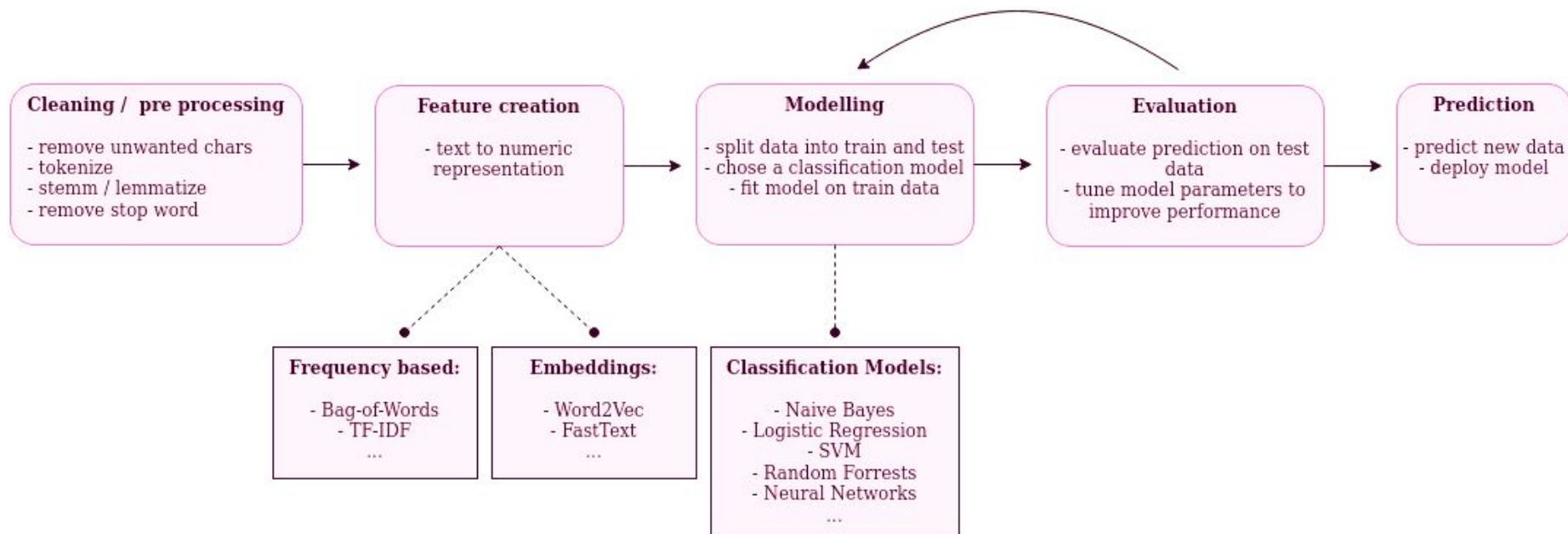
Feature Engineering & Column Vectorizer

“NLP is often applied for classifying text data. **Text classification** is the problem of assigning categories to text data according to its content. The most important part of text classification is **feature engineering**: the process of creating features for a machine learning model from raw text data.”

Column Vectorizer: The scikit-learn library offers easy-to-use tools to perform both tokenization and feature extraction of your text data.

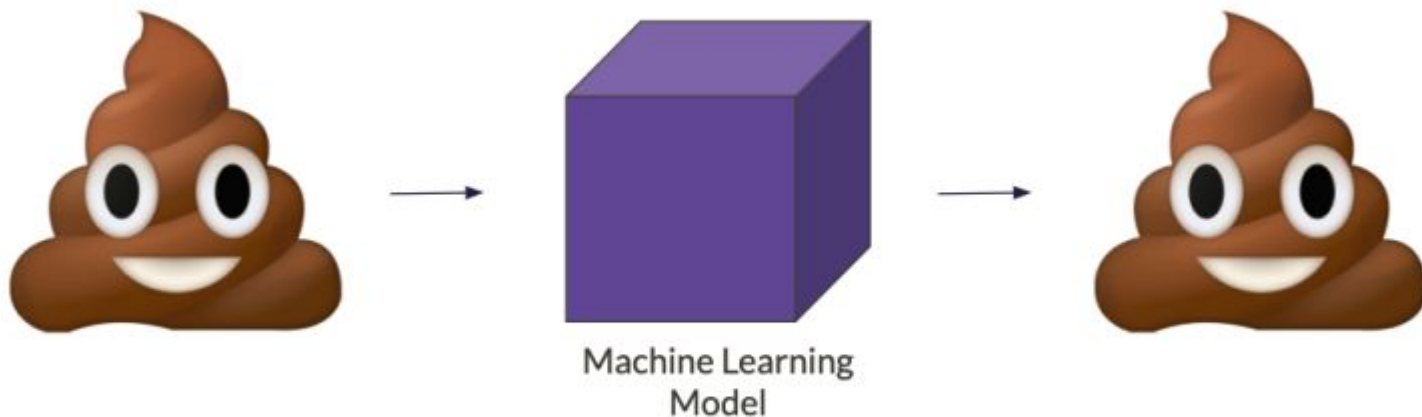


Typical NLP Workflow



The challenge

Data quality



- Garbage in garbage out
- Output quality depends on input quality

TayTweets the broken A.I.

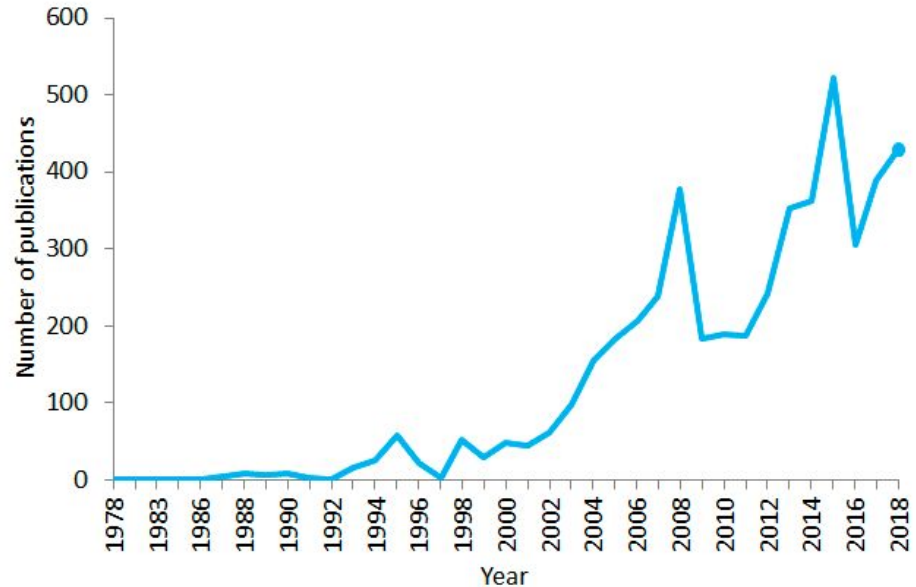


The Future

PubMed: 30 million biomedical citations/articles

-Big Data » need for structure and meaning » NLP.

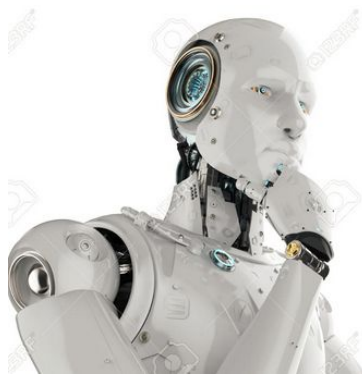
-Safe to conclude that as long as we will be dependent on Data, NLP is part of the process that allows for mass discoveries in an effortless manner.



Number of publications containing the sentence "natural language processing" in PubMed in the period 1978–2018. As of 2018, PubMed comprised more than 29 million citations for biomedical literature

Conclusion

- ❖ NLP is a subfield of ML that overlaps with other fields such as linguistics & C.S.
- ❖ NLP techniques are very powerful and used in a variety of different fields/scenarios:
- ❖ NLP aims at solving complex problems, and as a result big challenges and responsibilities comes as well. We should ensure quality control at every step.
- ❖ NLP is not just today's reality, but tomorrow's destiny as well.



Thank You!

References

TowardDS-<https://towardsdatascience.com/your-guide-to-natural-language-processing-nlp-48ea2511f6e1>,<https://towardsdatascience.com/text-analysis-feature-engineering-with-nlp-502d6ea9225d>

SAS-https://www.sas.com/en_us/insights/analytics/what-is-natural-language-processing-nlp.html

NLPFlow-<https://data-dive.com/german-nlp-binary-text-classification-of-reviews-part1>

MonkeyLearn- <https://monkeylearn.com/sentiment-analysis/>

YT-<https://www.youtube.com/watch?v=Lr4yi9onykg>

Book-<https://www.nltk.org/book/>