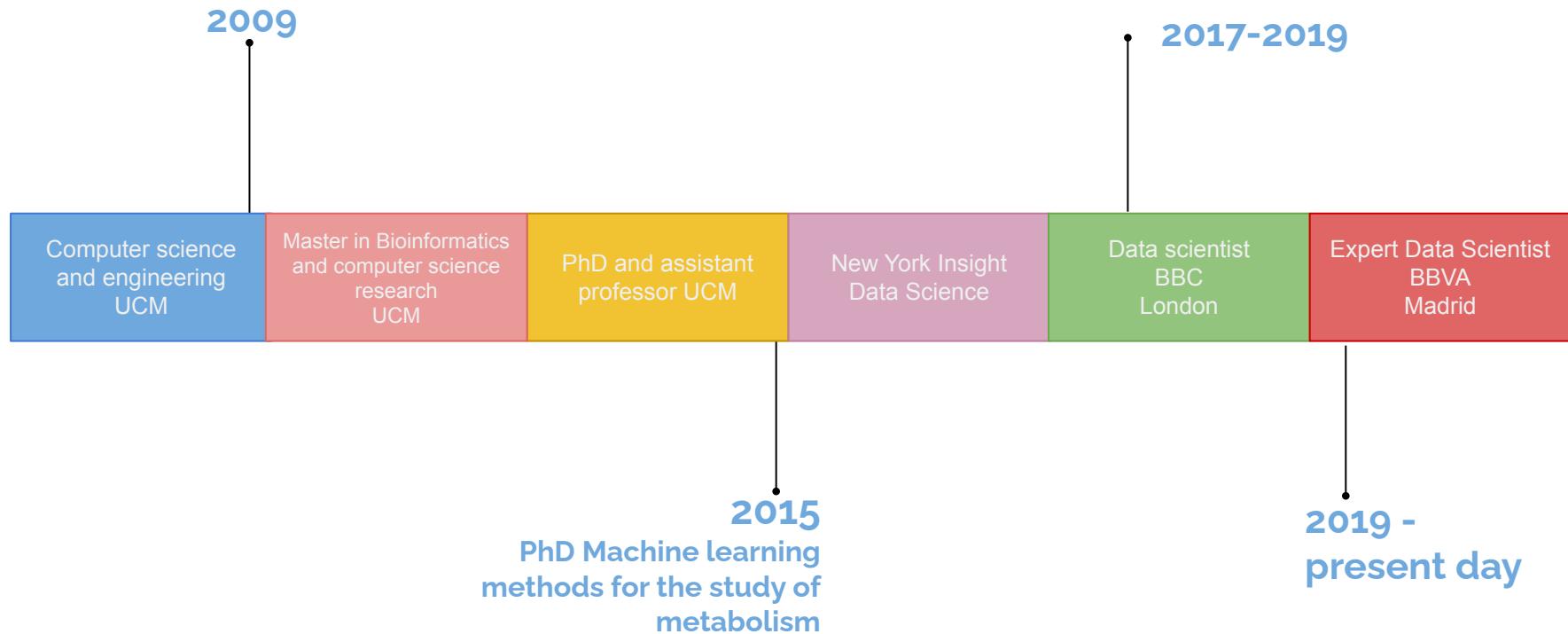

Introduction to Natural language processing

Session 1

Clara Higuera Cabañes, PhD
Barcelona Technology School, May 2021

About me



Class philosophy



Index

1. What is NLP?
2. NLP pipeline
 - a. Data Acquisition
 - b. Cleaning and preprocessing
 - c. Text representation / Feature engineering

Break (9:50, 10min)

3. Hands on / live showcase nlp in action (20-30min)
4. NLP pipeline (30 min)
 - a. Modeling
 - b. Evaluation of supervised learning problems

Break 10 mins

Quiz (15 min, 10-15 mins review collaboratively)

5. Hands on / live showcase nlp in action (30 min)
6. NLP real use case
 - a. Manual annotation

What is Natural Language Processing (NLP)?



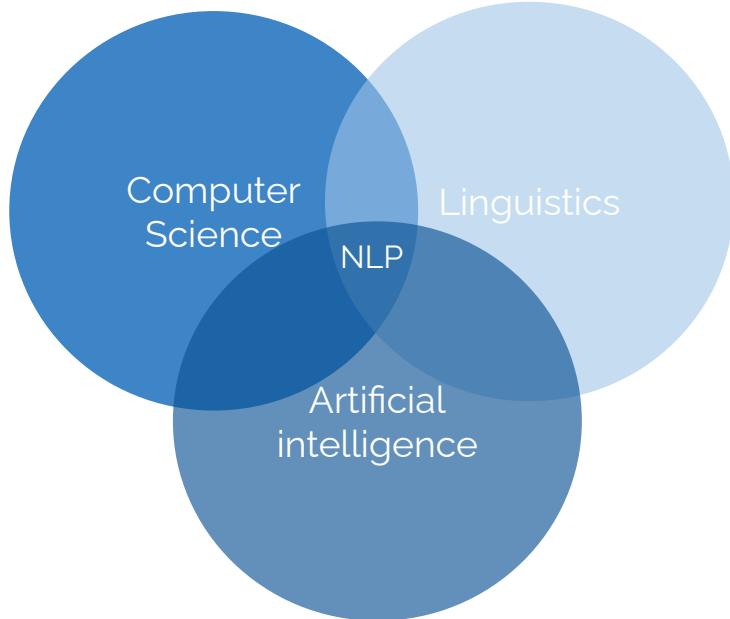
John: "How is the weather today?"

Digital assistant: "It is 37 degrees centigrade outside with no rain today."

John: "What does my schedule look like?"

Digital assistant: "You have a strategy meeting at 4 p.m. and an all-hands at 5:30 p.m. Based on today's traffic situation, it is recommended you leave for the office by 8:15 a.m."

What is Natural Language Processing?



It concerns building systems that can process and understand human language.

Sectors: retail, healthcare, finance, media, law, marketing, human resources, and many more.

Applications of NLP

Customer service

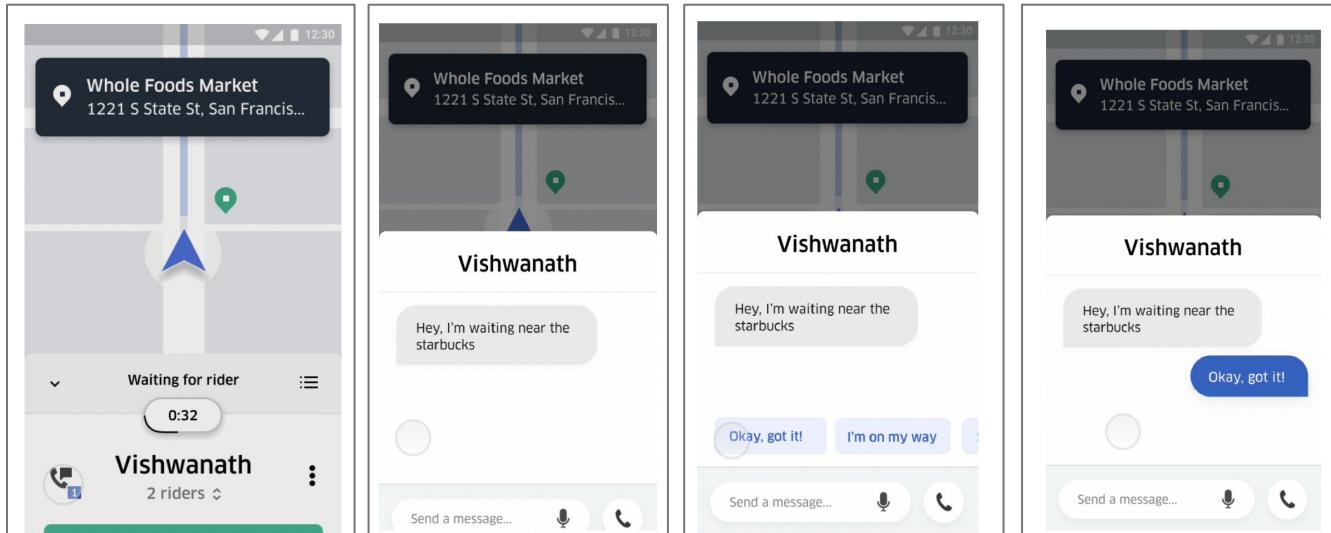
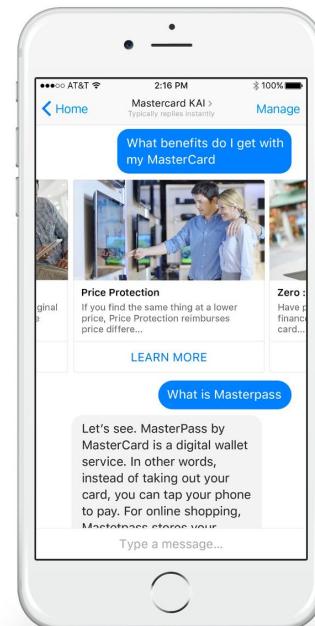
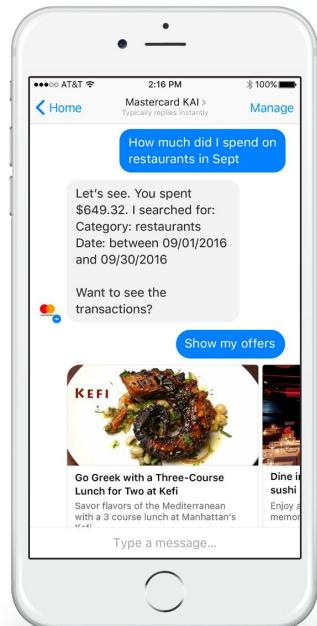
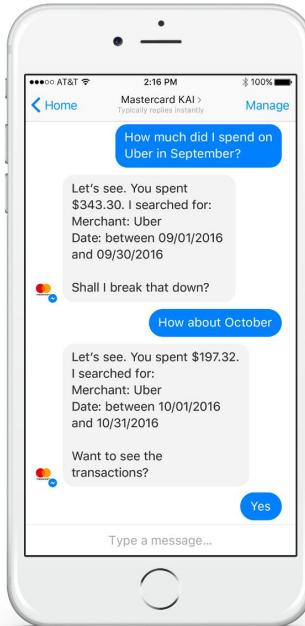


Figure 1: With one-click chat, driver-partners can more easily respond to rider messages.

Uber one-click-chat

Applications of NLP

Customer service



Applications of NLP

Question answering

who is the prime minister of spain

Aproximadamente 101.000.000 resultados (0,90 segundos)

Sugerencia: Buscar solo resultados en **español**. Puedes especificar tu idioma de búsqueda en Preferencias

España / Primer ministro

Pedro Sánchez

Desde 2018



Pedro Sánchez Pérez-Castejón es un político español, actual presidente del Gobierno de España. Es secretario general del Partido Socialista Obrero Español desde 2017, cargo que ya había desempeñado entre 2014 y 2016. [Wikipedia](#)

Nacimiento: 29 de febrero de 1972 (edad 49 años), Tetuán, Madrid

Estatura: 1,9 m

Cónyuge: María Begoña Gómez Fernández (m. 2006)

Hijas: Carlota Sánchez Gómez, Ainhoa Sánchez Gómez

Educación: Universidad Camilo José Cela (2012), MÁS

Padres: Magdalena Pérez-Castejón, Pedro Sánchez Sr.

who is the most famous spanish athlete

Aproximadamente 13.700.000 resultats (0,69 segons)

The Most Famous Spanish Athletes

- #1 Rafael Nadal. Born in Manacor in 1986, Nadal is currently ranked #1 by the Association of Tennis Professionals (ATP). ...
- #2 Pau Gasol. ...
- #3 Raúl González. ...
- #4 Fernando Alonso. ...
- #5 Miguel Indurain. ...
- #6 Mireia Belmonte. ...
- #7 Iker Casillas.

28 de nov. 2019

<https://www.gogoespana.com> › 2019/11/28 › famous-spa...

7 of the Most Famous Spanish Athletes of All Time

Altres persones també han preguntat

Who is the most famous Spanish person?

Pablo Picasso

1. Pablo Picasso. Pablo Picasso makes the number one spot in our list of **most famous Spanish people**. 18 d'abr. 2021

<https://gogoespana.com> › blog › top-15-most-famous-spa...

Top 15 most famous Spanish people - Go! Go! España

Cerca: Who is the most famous Spanish person?

Who is a famous Hispanic athlete?

Who is a Spanish athlete?

What is the most popular Spanish sport?

Who is the most famous Spanish singer?

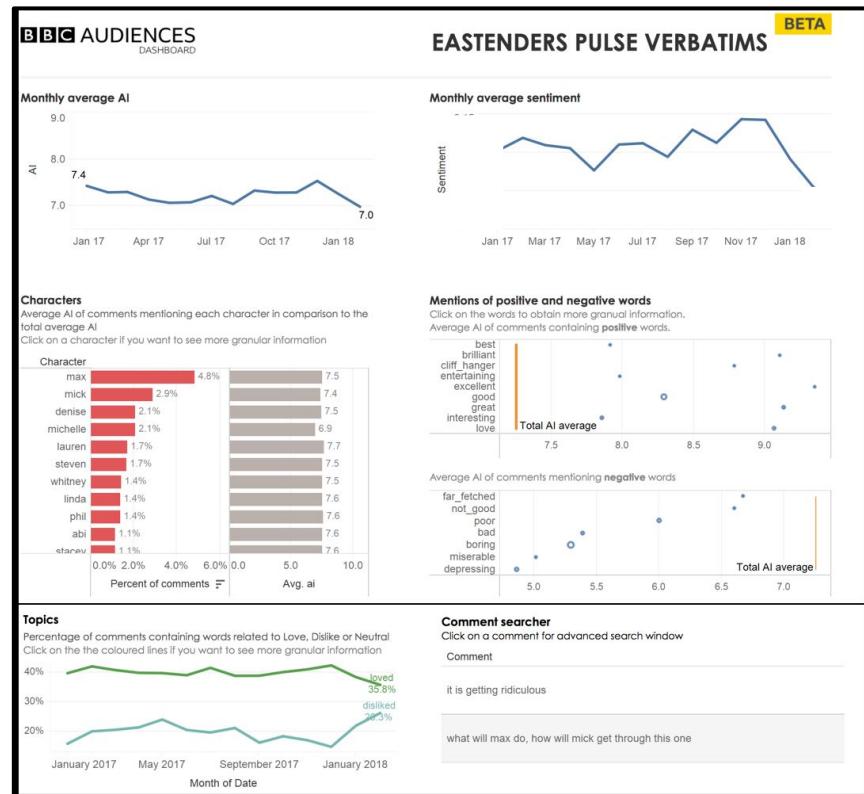
What are Spanish people called?

Applications of NLP

Media



B B C



Applications of NLP

Media



You might be interested as well in:

Content similarity recommenders

Applications of NLP

Media

You might be interested as well in:



Content similarity recommenders

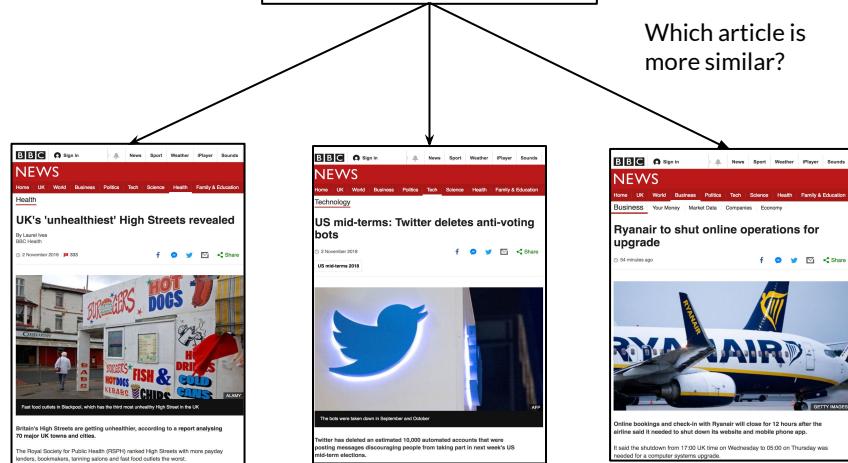
NHS prevention plan aims to boost life expectancy

Matt Hancock outlines new plans for an emphasis on illness prevention

People in England are being told to cut back on alcohol, sugar, salt and fat in a bid to boost the nation's life expectancy by five years.

Health Secretary Matt Hancock is setting out his long-term vision for the NHS on Monday - and the focus is on preventing illness.

Which article is more similar?



Applications of NLP

Medicine and social science

scientific reports

Explore content ▾ Journal information ▾ Publish with us ▾

nature > scientific reports > articles > article

Article | Open Access | Published: 09 May 2018

Identifying Suicide Ideation and Suicidal Attempts in a Psychiatric Clinical Research Database using Natural Language Processing

Andrea C. Fernandes , Rina Dutta, Sumithra Velupillai, Jyoti Sanyal, Robert Stewart & David Chandran

Scientific Reports 8, Article number: 7426 (2018) | Cite this article

8815 Accesses | 30 Citations | 14 Altmetric | Metrics

Abstract

Research into suicide prevention has been hampered by methodological limitations such as low sample size and recall bias. Recently, Natural Language Processing (NLP) strategies have been used with Electronic Health Records to increase information extraction from free text notes as well as structured fields concerning suicidality and this allows access to much larger cohorts than previously possible. This paper presents two novel NLP approaches – a rule-based approach to classify the presence of suicide ideation and a hybrid machine learning and rule-based approach to identify suicide attempts in a psychiatric clinical database. Good performance of the two classifiers in the evaluation study suggest they can be used to accurately detect mentions of suicide ideation and attempt within free-text documents in this psychiatric database. The novelty of the two approaches lies in the malleability of each classifier if a need to refine performance, or meet alternate classification requirements arises. The algorithms can also be adapted to fit infrastructures of other clinical datasets given sufficient clinical recording practice knowledge, without dependency on medical codes or additional data extraction of known risk factors to predict suicidal behaviour.

Applications of NLP

Reducing response times to citizen legal questions across Africa



barefootlaw

[Data Science for Social Good, Chicago University](#)

Understanding Text Data to Help Disadvantaged Families



Understanding Text Data to Help Disadvantaged Families

[DataKind UK](#)

<https://www.meetup.com/DataForGoodBCN/>

Improving forensics investigations

A growing number of government agencies are using NLP-based solutions to improve investigations in critical areas such as law enforcement, defense, and intelligence. The DoD's DEFT program referenced above uses NLP to uncover connections implicit in large text documents. Its objective is to improve the efficiency of defense analysts who investigate multiple documents to detect anomalies and causal relationships.³⁶

The European Union's Horizon 2020 program launched an initiative called RED (Real-time Early Detection) Alert, aimed at countering terrorism by using NLP to monitor and analyze social media conversations. RED Alert is designed to provide early alerts of potential propaganda and signs of warfare by identifying online content posted by extremists. To comply with the General Data Protection Regulation (GDPR), this analysis uses homomorphic encryption, a method that allows mathematical operations to be performed on encrypted text without disturbing the original encryption.³⁷

Enhancing policy analysis

The World Bank's Poverty and Equity Global Practice Group used LDA topic modeling to measure changes in policy priorities by examining presidential speeches in 10 Latin American countries and Spain from 1819 to 2016. Using LDA, the authors could identify the main topics for each document and indicate the variation in their significance across countries and over time. In Peru, for instance, topics on infrastructure and public services diminished in importance over time. With the help of topic modeling, the authors were able to establish, for each nation, a negative correlation between policy volatility and long-term growth.⁴¹

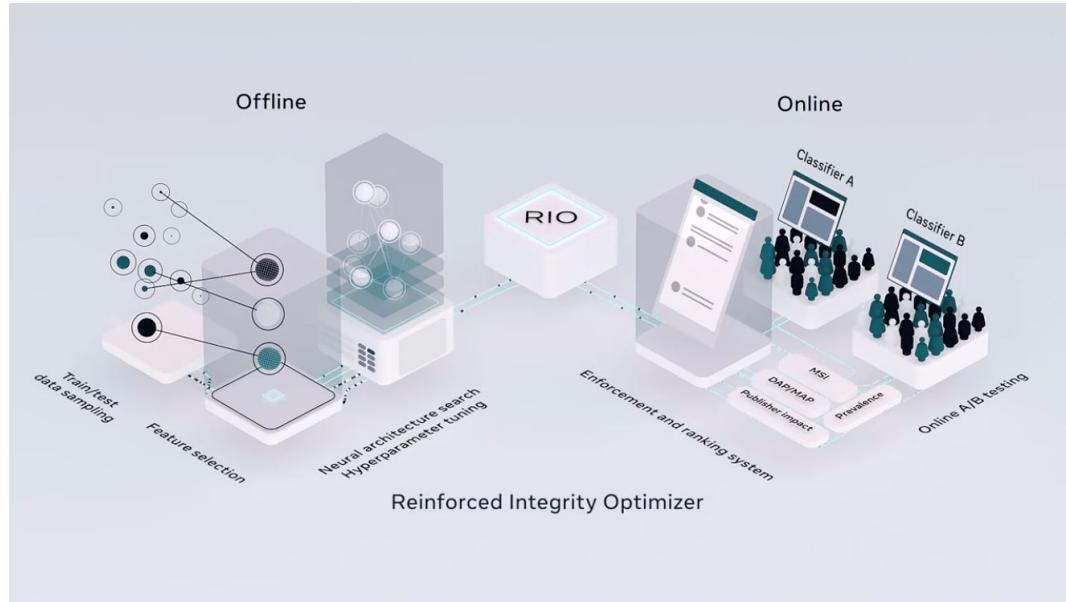
Improving predictions to aid decision-making

One of the most striking characteristics of NLP is its ability to facilitate better predictions, which can help agencies design preemptive measures. The police department of Durham, North Carolina, uses NLP in crimefighting by enabling the police to observe patterns and interrelations in criminal activities and identify pockets with a high incidence of crime, thus allowing for quicker interventions. This contributed to a 39 percent drop in violent crime in Durham from 2007 to 2014.³⁸

NLP also is being used to combat child trafficking. About 75 percent of child trafficking involves online advertisements. DARPA, in collaboration with commercial technology experts, has developed a platform that monitors and draws connections among the dubious content of online advertisements. Virginia's Fairfax County Police Department and New Orleans's Homeland Security investigations both use this advanced software to identify high-risk web advertisements and detect code words used by traffickers.³⁹

Applications of NLP

Hate Speech detection



[Facebook](#)

NLP is developing quickly

- **Widely available and easy-to-use NLP tools**, techniques, and APIs are now all-pervading in the industry. There has never been a better time to build quick NLP solutions.
- **Development of more interpretable and generalized approaches** has improved the baseline performance for even complex NLP tasks, such as open-domain conversational tasks and question answering, which were not practically feasible before.
- More and **more organizations**, including Google, Microsoft, and Amazon, **are investing heavily in more interactive consumer products, where language** is used as the primary medium of communication.
- **Increased availability of useful open source datasets**, along with standard benchmarks on them, has acted as a catalyst in this revolution, as opposed to being impeded by proprietary datasets only available to limited organizations and individuals.
- **The viability of NLP has moved beyond English or other major languages.** Datasets and language-specific models are being created for the less-frequently digitized languages too. A fruitful product that came out this effort was a near-perfect automatic machine translation tool available to all individuals with a smartphone.
- The **amount of data available plus the advances in deep learning** is 'opening' a new era in NLP

NLP capabilities

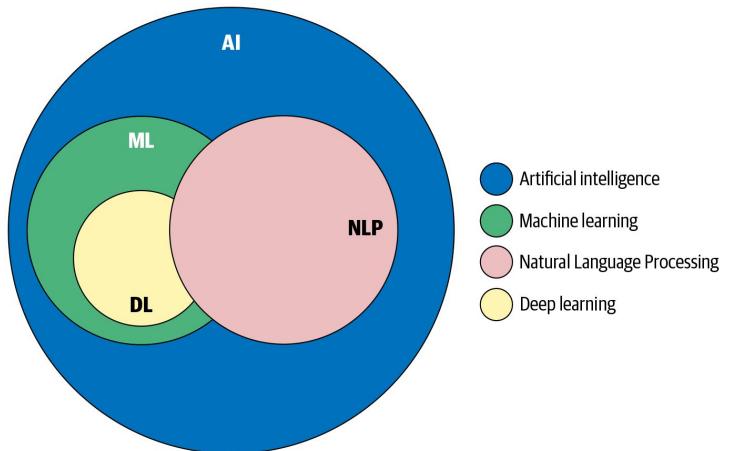
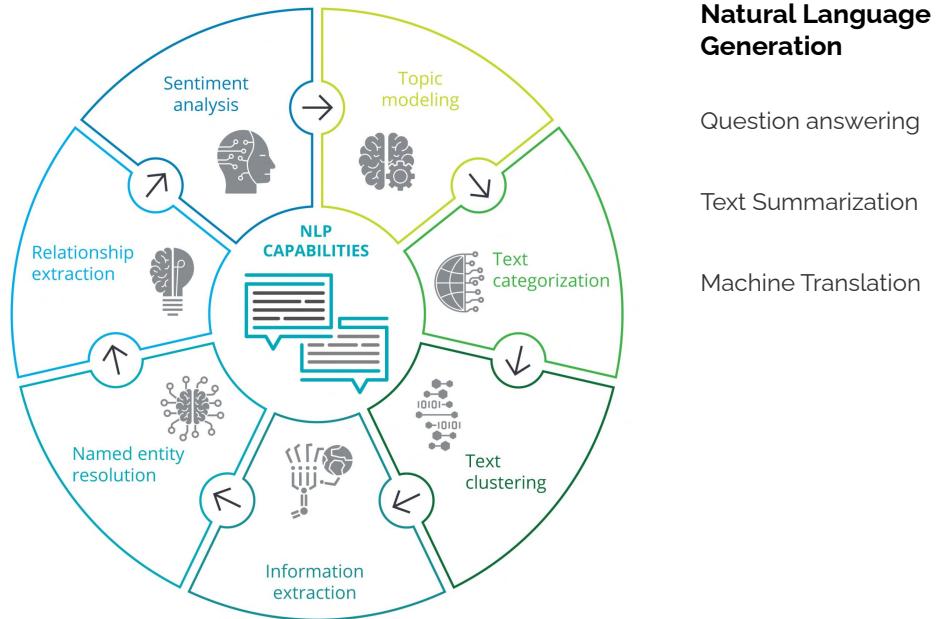


FIGURE 2
Key NLP capabilities



Source: Deloitte analysis.

Deloitte Insights | deloitte.com/insights

Natural Language Generation

Question answering

Text Summarization

Machine Translation

NLP Tools and Libraries

spaCy

spaCy Out now: spaCy v3.0

GET STARTED Installation

- Quickstart
- Instructions
- Troubleshooting
- Changelog

Models & Languages

- Facts & Figures
- spaCy 101
- New in v3.0

GUIDES

- Linguistic Features
- Rule-based Matching
- Processing Pipelines
- Embeddings & Transformers NEW
- Training Models NEW
- Layers & Model Architectures NEW
- spaCy Projects NEW
- Saving & Loading
- Visualizers

RESOURCES

- Project Templates
- v2.x Documentation

Install spaCy

macOS / OSX Windows Linux

Package manager

- pip conda from source

Hardware

- CPU GPU

Configuration

- virtual env ? train models ?

Trained pipelines

- Chinese Danish Dutch English French German Greek Italian Japanese Lithuanian Multi-language Norwegian Bokmål Polish Portuguese Romanian Russian Spanish

Select pipeline for

- efficiency ? accuracy ?

```
$ pip install -U pip setuptools wheel
$ pip install -U spacy
$ python -m spacy download en_core_web_sm
```

<https://spacy.io/usage>



NLTK 3.6.2 documentation

[PREVIOUS](#) | [NEXT](#) | [MODULES](#) | [INDEX](#)

Installing NLTK

NLTK requires Python versions 3.5, 3.6, 3.7, 3.8, or 3.9

For Windows users, it is strongly recommended that you go through this guide to install Python 3 successfully <https://docs.python-guide.org/starting/install3/win/#install3-windows>

Setting up a Python Environment (Mac/Unix/Windows)

Please go through this guide to learn how to manage your virtual environment managers before you install NLTK, <https://docs.python-guide.org/dev/virtualenvs/>

Alternatively, you can use the Anaconda distribution installer that comes "batteries included" <https://www.anaconda.com/distribution/>

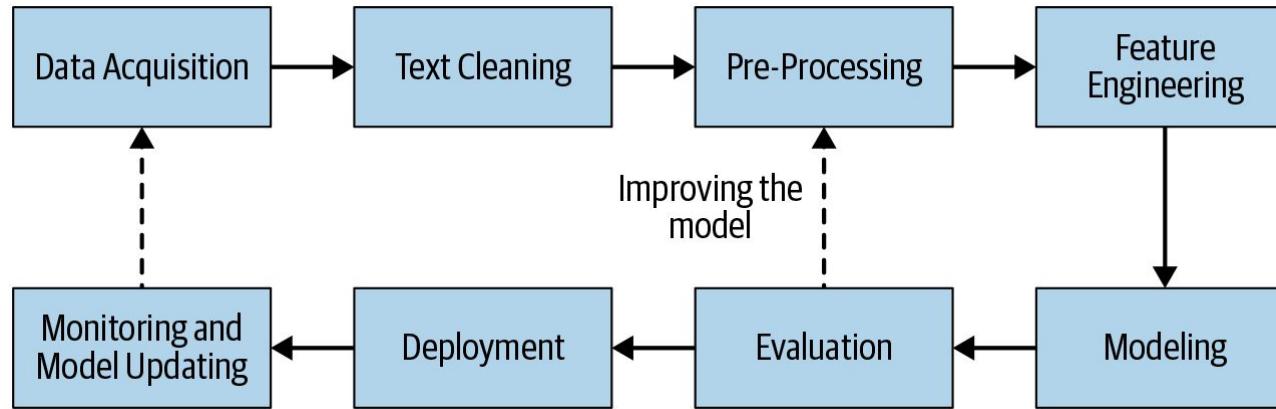
Mac/Unix

1. Install NLTK: run `pip install --user -U nltk`
2. Install Numpy (optional): run `pip install --user -U numpy`
3. Test installation: run `python` then type `import nltk`

For older versions of Python it might be necessary to install setuptools (see <http://pypi.python.org/pypi/setup-tools>) and to install pip (`sudo easy_install pip`).

<https://www.nltk.org/install.html>

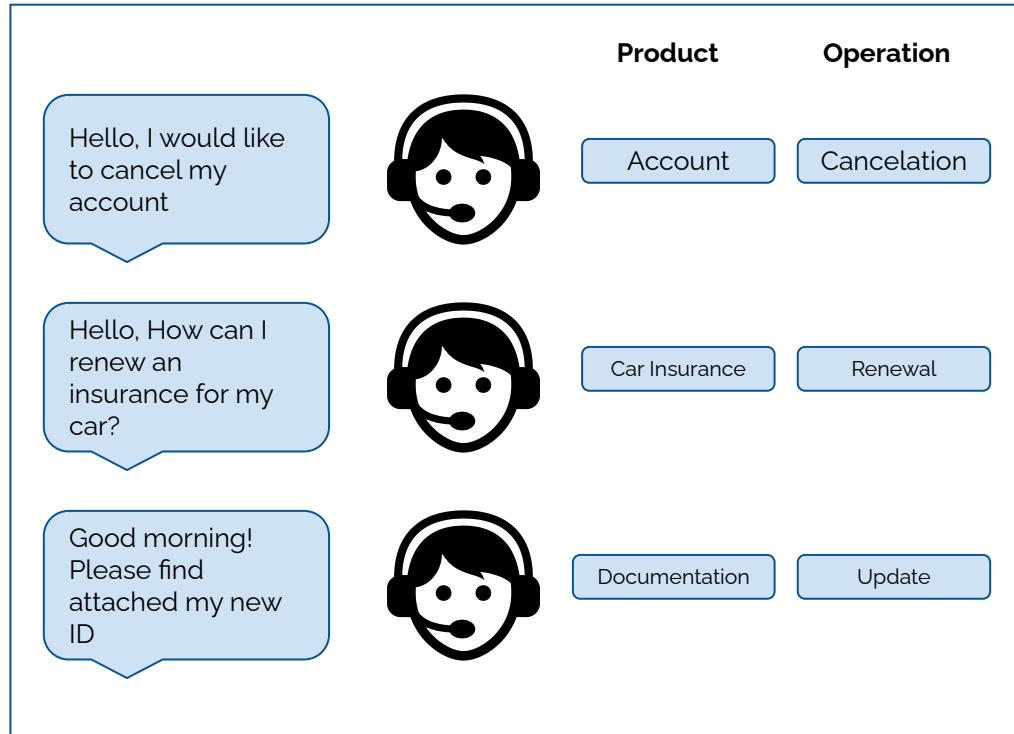
NLP Pipeline



NLP Pipeline

Data acquisition

- Use a public dataset
- Scrape data
- Product intervention / sensorization
- Data Augmentation -> refer to paper from EMNLP!
 - Back translation
 - Active learning



NLP Pipeline

Text preprocessing

Depends highly on the use case

- Language detection (polyglot)
- Text normalization
 - Convert all text to lowercase or uppercase
 - Convert digits to text
 - Normalize date format
- Remove punctuation, digits
- **Remove stop words**
- Sentence segmentation
- Stemming / lemmatization
- Word tokenization
- Bigram / trigram..
- Advanced processing
 - POS tagging, entity recognition

Sentence segmentation

```
text = u"This is first sentence. Second sentence. Third sentence."  
text_sentences = nlp(text)  
for sentence in text_sentences.sents:  
    print(sentence.text)
```



This is first sentence.
Second sentence.
Third sentence.

NLP Pipeline

Text preprocessing

Depends highly on the use case

- Language detection (polyglot)
- Text normalization
 - Convert all text to lowercase or uppercase
 - Convert digits to text
 - Normalize date format
- Remove punctuation, digits
- **Remove stop words**
- Sentence segmentation
- Stemming / lemmatization
- Word tokenization
- Bigram / trigram..
- Advanced processing
 - POS tagging, entity recognition

Stop words

“ Generally, the most common words used in a text are “the”, “is”, “in”, “for”, “where”, “when”, “to”, “at” etc.

NLP Pipeline

Text preprocessing

Depends highly on the use case

- Language detection (polyglot)
- Text normalization
 - Convert all text to lowercase or uppercase
 - Convert digits to text
 - Normalize date format
- Remove punctuation, digits
- Remove stop words
- Sentence segmentation
- **Stemming / lemmatization**
- Word tokenization
- Bigram / trigram..
- Advanced processing
 - POS tagging, entity recognition

Stemming

adjustable -> adjust
formality -> formaliti
formaliti -> formal
airliner -> airlin

Lemmatization

was -> (to) be
better -> good
meeting -> meeting

```
nlp = spacy.load('en_core_web_sm')
doc = nlp(u'I am Clara and this is the NLP class')
#for token in doc:
#    print(token.text, token.lemma_, token.is_stop)

for token in doc:
    print(token.text, token.lemma_)
```



```
(u'I', u'-PRON-', True)
(u'am', u'be', True)
(u'Clara', u'Clara', False)
(u'and', u'and', True)
(u'this', u'this', True)
(u'is', u'be', True)
(u'the', u'the', True)
(u'NLP', u'NLP', False)
(u'class', u'class', False)
```

NLP Pipeline

Text preprocessing

Depends highly on the use case

- Language detection (polyglot)
- Text normalization
 - Convert all text to lowercase or uppercase
 - Convert digits to text
 - Normalize date format
- Remove punctuation, digits
- Remove stop words
- Sentence segmentation
- Stemming / lemmatization
- Word tokenization
- Bigram / trigram..
- **Advanced processing**
 - POS tagging, entity recognition

Input

Chaplin wrote, directed, and composed the music for most of his films.

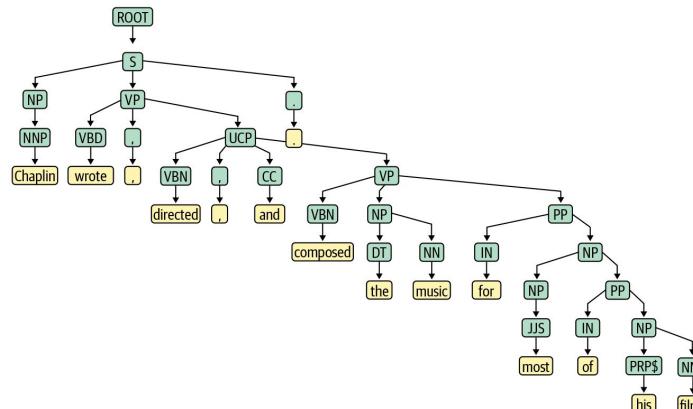
Tokenization with Lemmatization

Chaplin write direct and compose the music for most of he film
Chaplin wrote, directed, and composed the music for most of his films.

POS Tagging

NNP VBD . VBD CC VBN DT NN IN JJS IN PRPS NNS
Chaplin wrote, directed, and composed the music for most of his films.

Parse Tree



NLP Pipeline

Text preprocessing

Depends highly on the use case

- Language detection (polyglot)
- Text normalization
 - Convert all text to lowercase or uppercase
 - Convert digits to text
 - Normalize date format
- Remove punctuation, digits
- Remove stop words
- Sentence segmentation
- Stemming / lemmatization
- Word tokenization
- Bigram / trigram..
- **Advanced processing**
 - POS tagging, entity recognition

```
nlp = spacy.load('en_core_web_sm')

doc = nlp(u'I am Clara and this is the NLP class')

for token in doc:
    print(token.text, token.lemma_, token.is_stop, token.pos_)
```

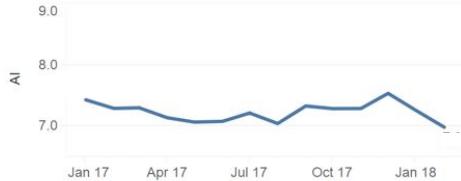


```
(u'I', u'-PRON-', True, u'PRON')
(u'am', u'be', True, u'VERB')
(u'Clara', u'Clara', False, u'PROPN')
(u'and', u'and', True, u'CCONJ')
(u'this', u'this', True, u'DET')
(u'is', u'be', True, u'VERB')
(u'the', u'the', True, u'DET')
(u'NLP', u'NLP', False, u'PROPN')
(u'class', u'class', False, u'NOUN')
```

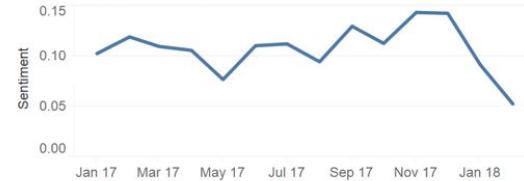
EASTENDERS PULSE VERBATIMS

BETA

Monthly average AI

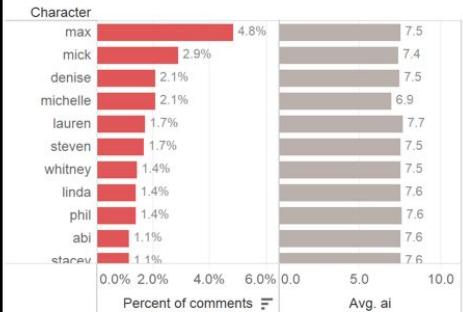


Monthly average sentiment



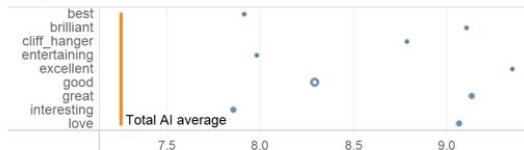
Characters

Average AI of comments mentioning each character in comparison to the total average AI
Click on a character if you want to see more granular information

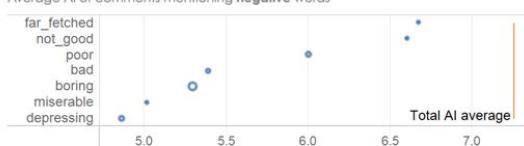


Mentions of positive and negative words

Click on the words to obtain more granular information.
Average AI of comments containing positive words.

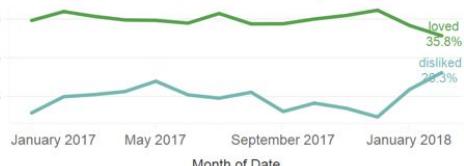


Average AI of comments mentioning negative words



Topics

Percentage of comments containing words related to Love, Dislike or Neutral
Click on the coloured lines if you want to see more granular information



Comment searcher

Click on a comment for advanced search window

Comment

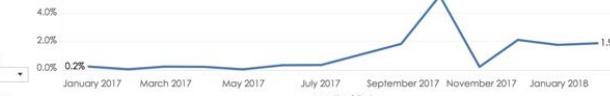
it is getting ridiculous

what will max do, how will mick get through this one

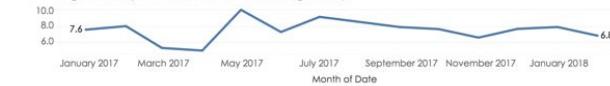
BETA

Home button

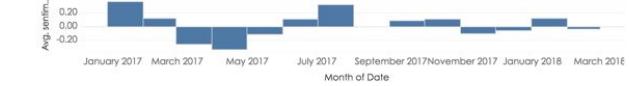
Percentage of comments mentioning stacey



Average monthly AI of comments mentioning stacey



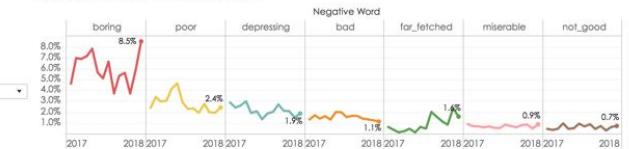
Average monthly sentiment of comments mentioning stacey



BETA

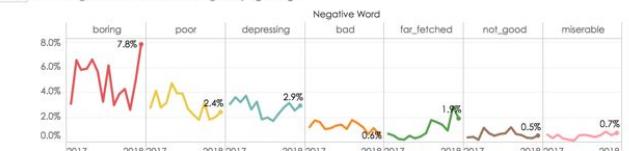
Home button

Total percentage of comments mentioning All

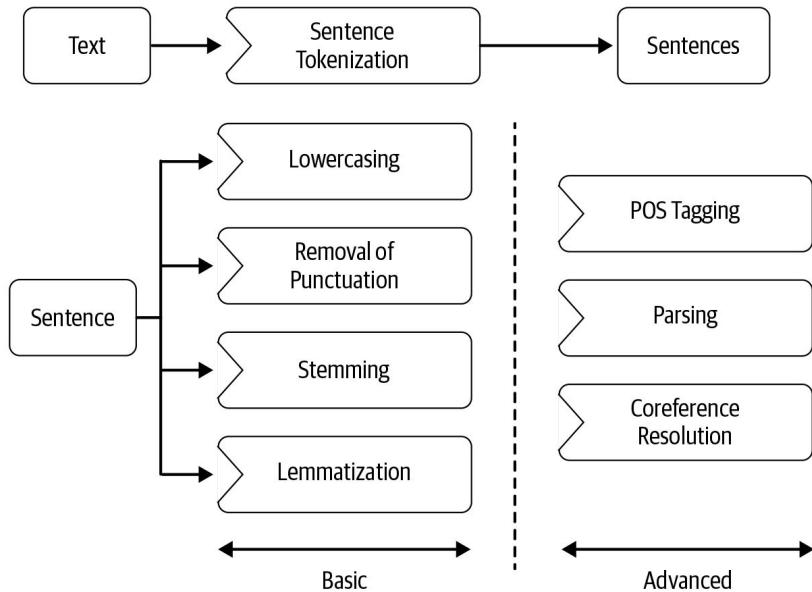


Filter by age range

Percentage of comments mentioning All by age range

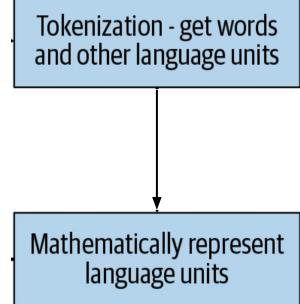
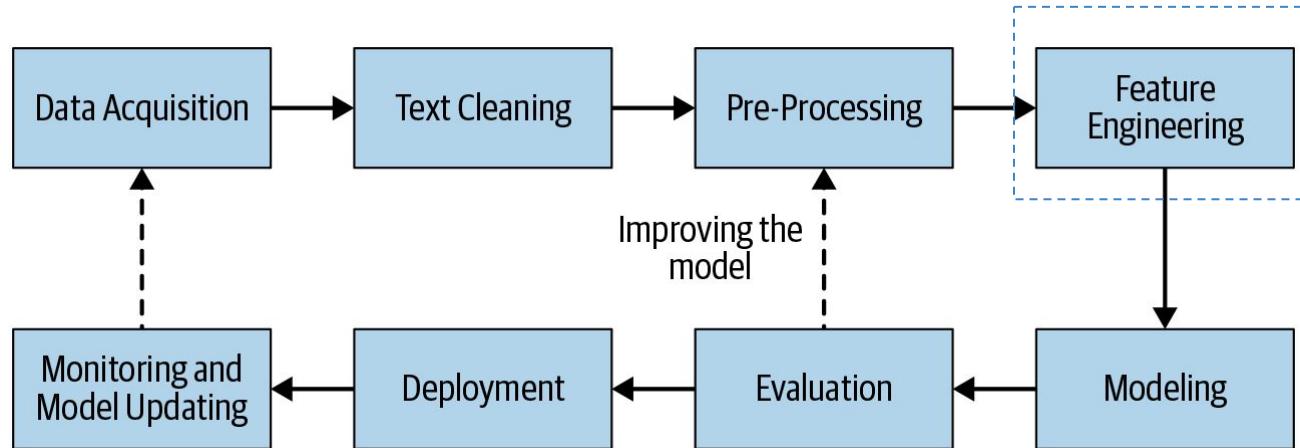


NLP Pipeline



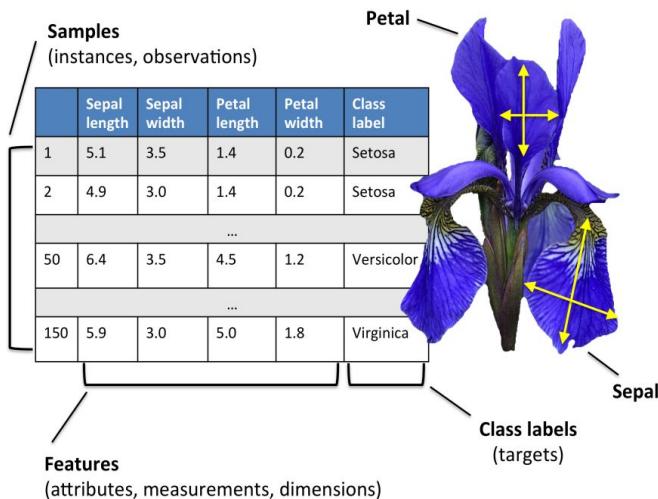
NLP Pipeline

Feature engineering or text representation

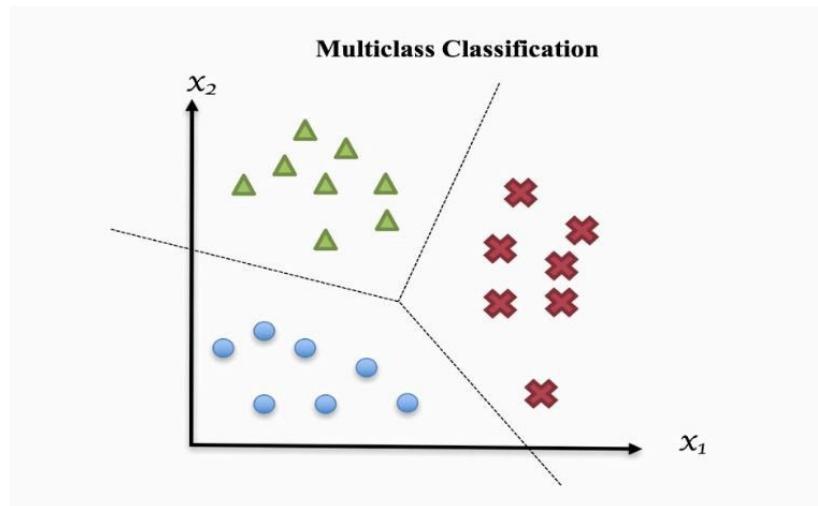


NLP Pipeline

Feature engineering or text representation



Classical input dataset for a ML algorithm



How do we train a model with text?
How do we translate text into a mathematical representation?

Bag of words approach (BoW)

John likes to watch movies. Mary likes movies too.
Mary also likes to watch football games.

```
{"John":1,"likes":3,"to":2,"watch":2,"movies":2,"Mary":2,"too":1,"also":1,"football":1,"games":1}
```

	the	red	dog	cat	eats	food
1. the red dog	1	1	1	0	0	0
2. cat eats dog	0	0	1	1	1	0
3. dog eats food	0	0	1	0	1	1
4. red cat eats	0	1	0	1	1	0

→ Vocabulary

```
from sklearn.feature_extraction.text import  
CountVectorizer  
count_vect = CountVectorizer(binary=True)  
bow_rep_bin = count_vect.fit_transform(processed_docs)  
temp = count_vect.transform(["the red dog"])  
  
print("Bow representation for 'the red dog':",  
temp.toarray())
```

Tokenization - get words and other language units

Mathematically represent language units

Revisar
binary=True

Hands on exercise...

How do we train a model with text?

How do we translate text into a mathematical representation?

Bag of words approach (BoW) and **TF-IDF**

$$w_{i,j} = tf_{i,j} \times \log \left(\frac{N}{df_i} \right)$$

$tf_{i,j}$ = number of occurrences of i in j
 df_i = number of documents containing i
 N = total number of documents

```
from sklearn.feature_extraction.text import TfidfVectorizer
tfidf = TfidfVectorizer()
bow_rep_tfidf = tfidf.fit_transform(processed_docs)
print(tfidf.idf_) #IDF for all words in the vocabulary
print(tfidf.get_feature_names()) #All words in the vocabulary.

temp = tfidf.transform(["dog and man are friends"])

print("Tfidf representation for 'dog and man are friends'\n",
      temp.toarray())
```

- Vectors instead of binary values have real values
- Penalises very frequent words

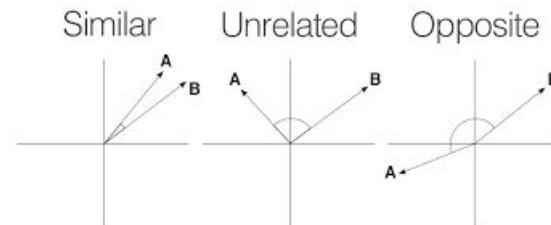
How do we train a model with text?
How do we translate text into a mathematical representation?

Bag of words approach (BoW) and **TF-IDF**

Vectorizing texts allows to make mathematical calculations with it

$$\text{similarity}(A, B) = \frac{A \cdot B}{\|A\| \times \|B\|} = \frac{\sum_{i=1}^n A_i \times B_i}{\sqrt{\sum_{i=1}^n A_i^2} \times \sqrt{\sum_{i=1}^n B_i^2}}$$

Cosine similarity

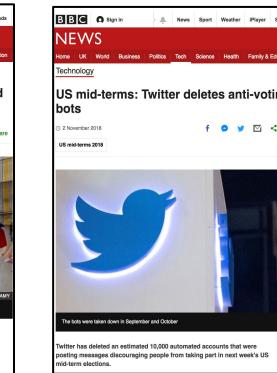


How do we train a model with text?

How do we translate text into a mathematical representation?

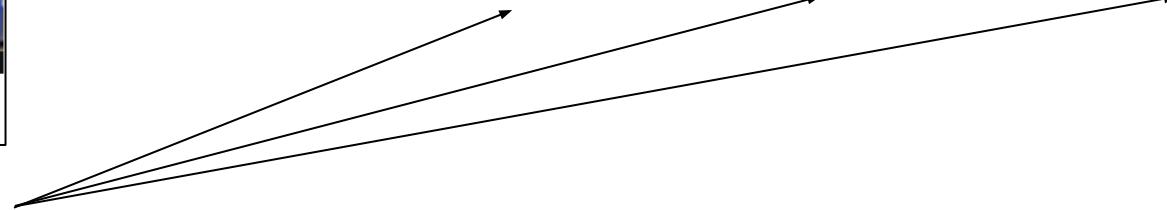
Bag of words approach (BoW) and **TF-IDF**

Vectorizing texts allows to make mathematical calculations with it

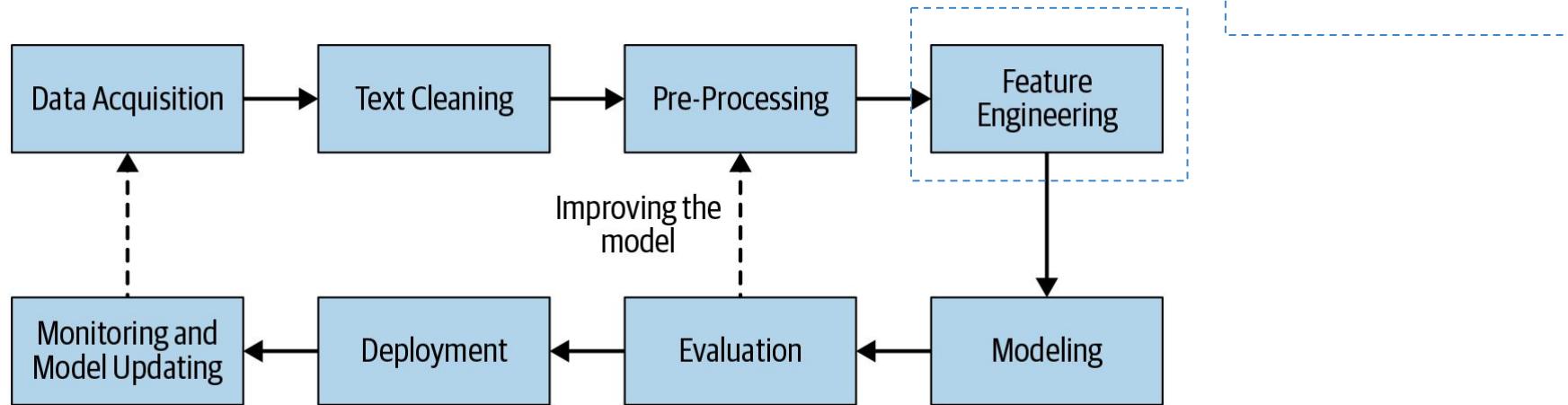


Article read

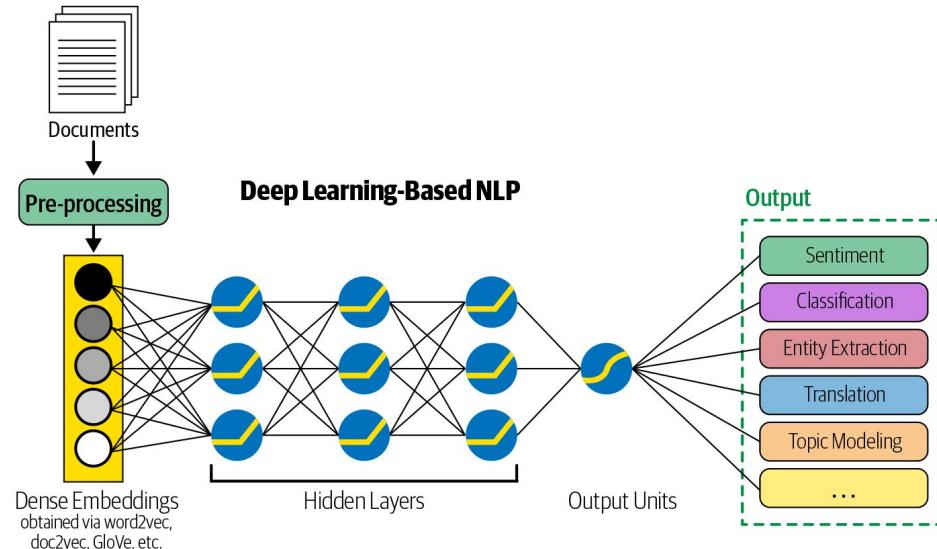
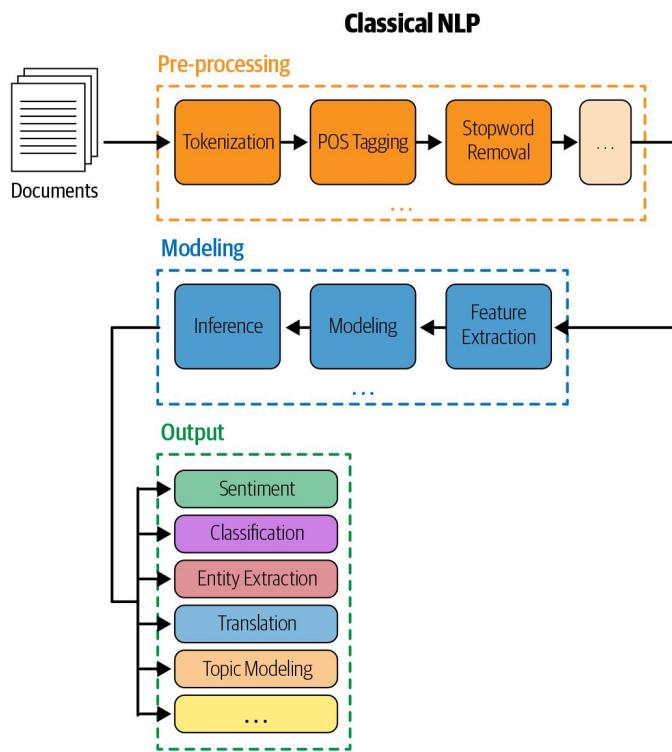
$(a_1, a_2, a_3, a_4, \dots, a_n)$



NLP Pipeline



NLP Pipeline

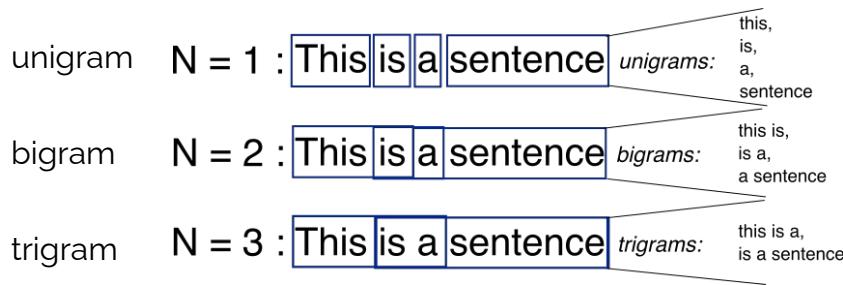


How do we train a model with text? How do we translate text into a mathematical representation?

Bag of words approach (BoW)

- + Transforms text into vectors
- + Simplicity
- Ignores order and context
- Size of the vector increases with the size of the vocabulary
- Does not capture the similarity between words that mean the same (i.e. run and ran)

We can also use bigrams, trigrams,...



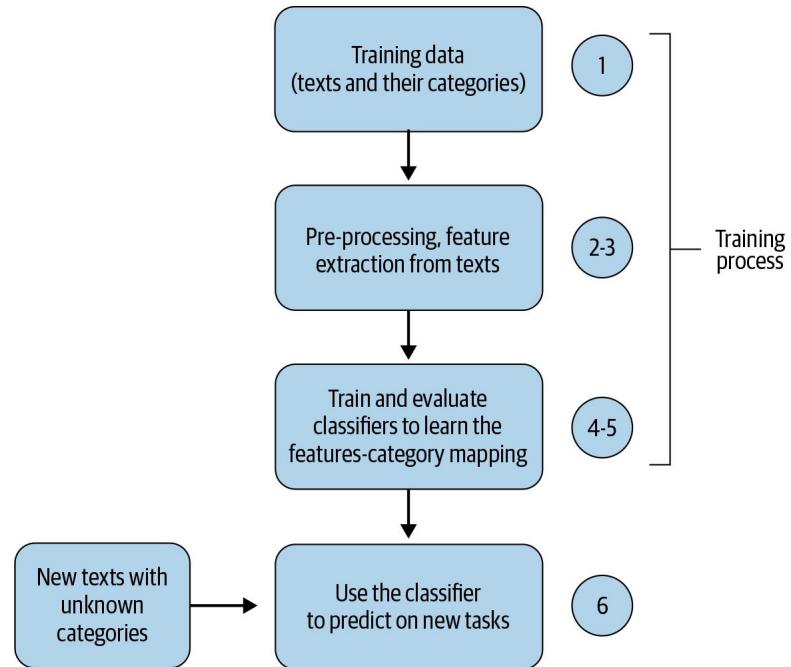
#n-gram vectorization example with count vectorizer and uni, bi, trigrams

```
count_vect = CountVectorizer(ngram_range=(1,3))
```


NLP pipeline

Modeling

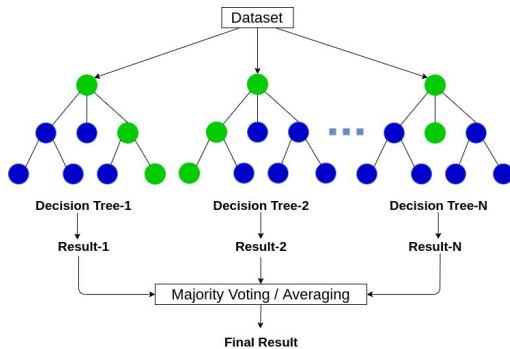
- **Types of problems**
 - Supervised learning problems
 - Unsupervised learning problems
 - Language generation problems
- Evaluation (metrics)
 - Supervised learning problems
 - Unsupervised learning problems
 - Generated language problems



NLP pipeline

Modeling

- Types of problems
 - Supervised learning problems - **classification**

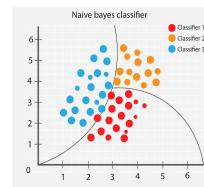


Random Forest

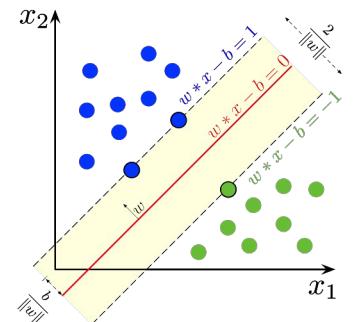
$$P(A|B) = \frac{P(B|A) P(A)}{P(B)}$$

using Bayesian probability terminology, the above equation can be written as

$$\text{Posterior} = \frac{\text{prior} \times \text{likelihood}}{\text{evidence}}$$



Naive Bayes



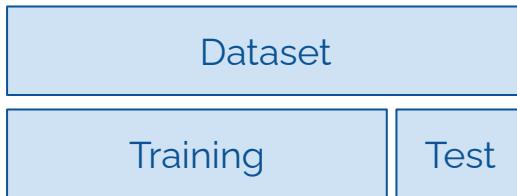
Support vector machines

Same models as in non-text problems

NLP Pipeline

Modeling

- Evaluation
 - Supervised learning problems



Regression	Classification
<ul style="list-style-type: none">• Mean Absolute Error (MAE)• Root Mean Squared Error (RMSE)• R-Squared and Adjusted R-Squared	<ul style="list-style-type: none">• Recall• Precision• F1-Score• Accuracy• Area Under the Curve (AUC)

Hands on exercise...

Quiz (15 min, 10-15 mins review collaboratively)

Introduction to Natural language processing

Session 2

Clara Higuera Cabañes, PhD
Barcelona Technology School, May 2021

Index

1. Pending concepts from session 1 + unsupervised learning NLP (1h 30min)

2. Unsupervised learning in NLP (30 min)
 - a. Clustering
 - b. Topic modeling
 - c. Distance metrics
3. LDA applied
 - a. Recommender system for BBC News with topic modeling
 - b. Reading profile of users for segmentation
 - c. Political knowledge gaps

Break 10 mins

4. Hands on / live showcase nlp in action - LDA notebook (30 min)

5. Deep learning (1h 30min)

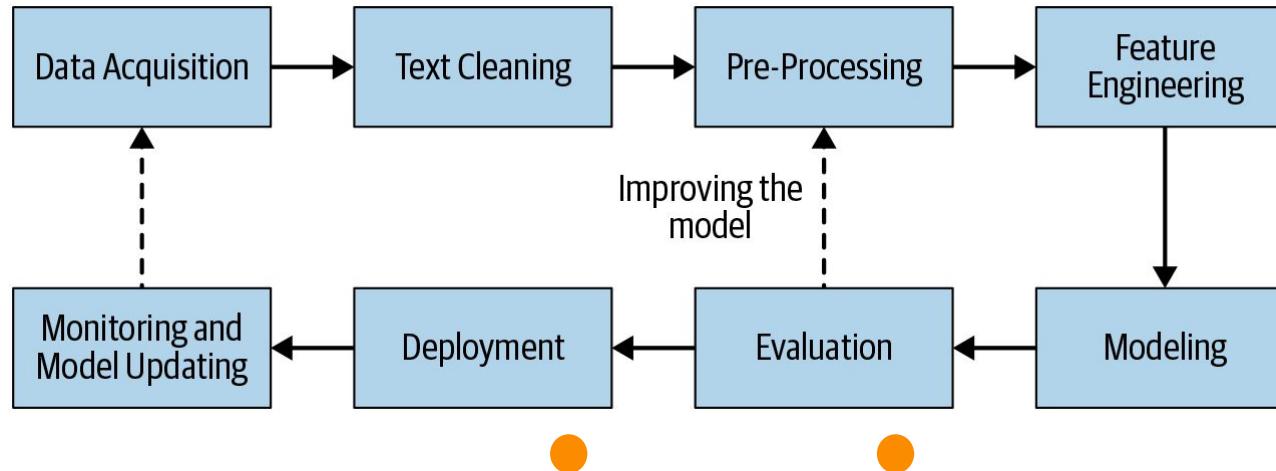
6. Word2vec - Embeddings (30 min)
 - a. Notebook (15 min)

Break 10 min

7. Text classification with deep learning (15-20 min)
8. Natural language generation (20 min)

Quiz (20 min + 10 review)

Previously on NLP pipeline ...

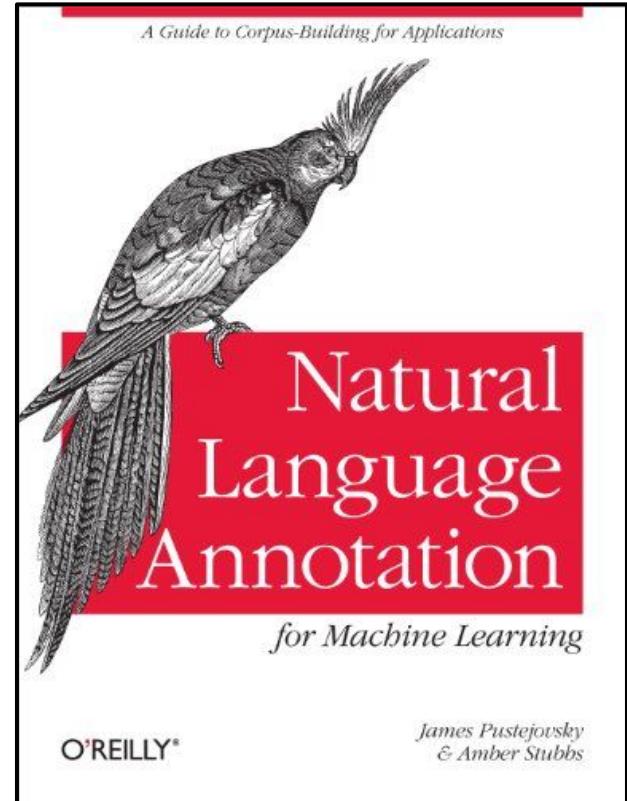


● Concepts pending
from session 1

NLP Pipeline

Modeling

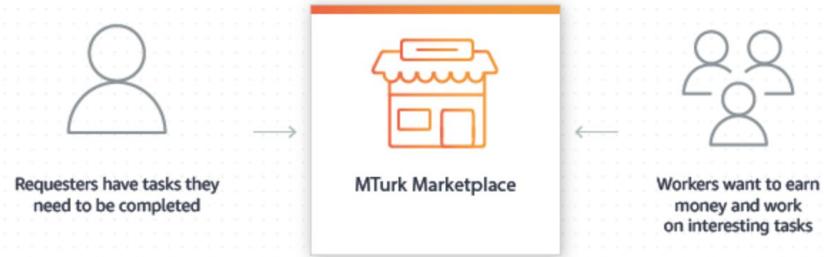
- Manual annotation
 - Needed when:
 - Lack of data
 - NLP models validation



NLP Pipeline

Modeling

- **Manual annotation**
 - **Needed when:**
 - Lack of data
 - NLP models validation
 - **DIY / outsource**

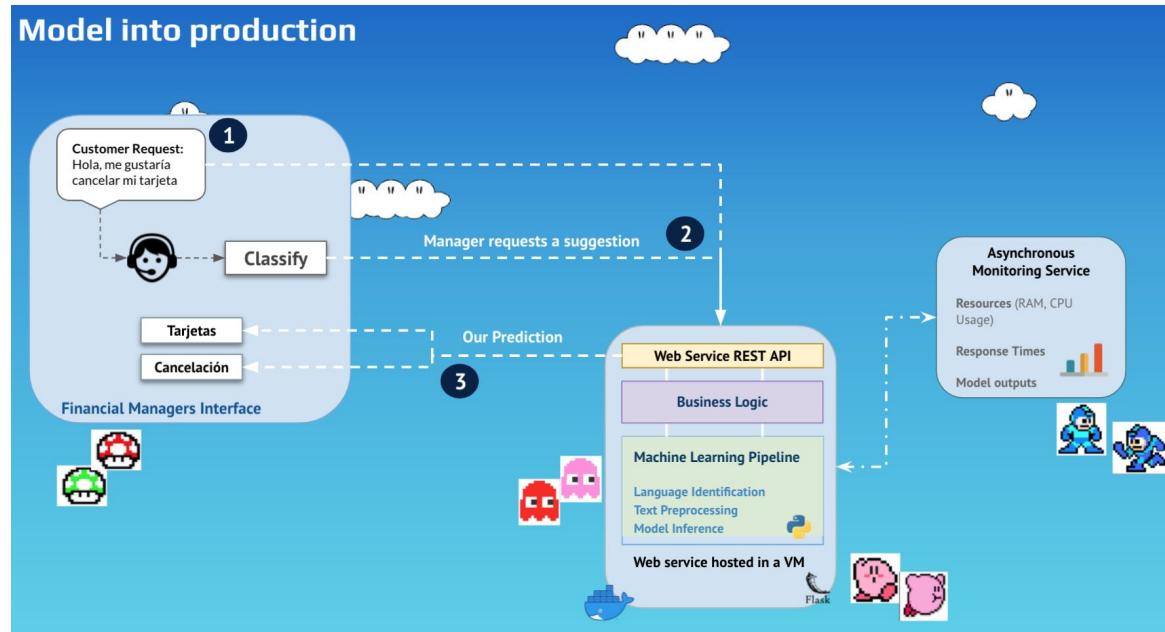


[Amazon Mechanical Turk](#)

NLP Pipeline

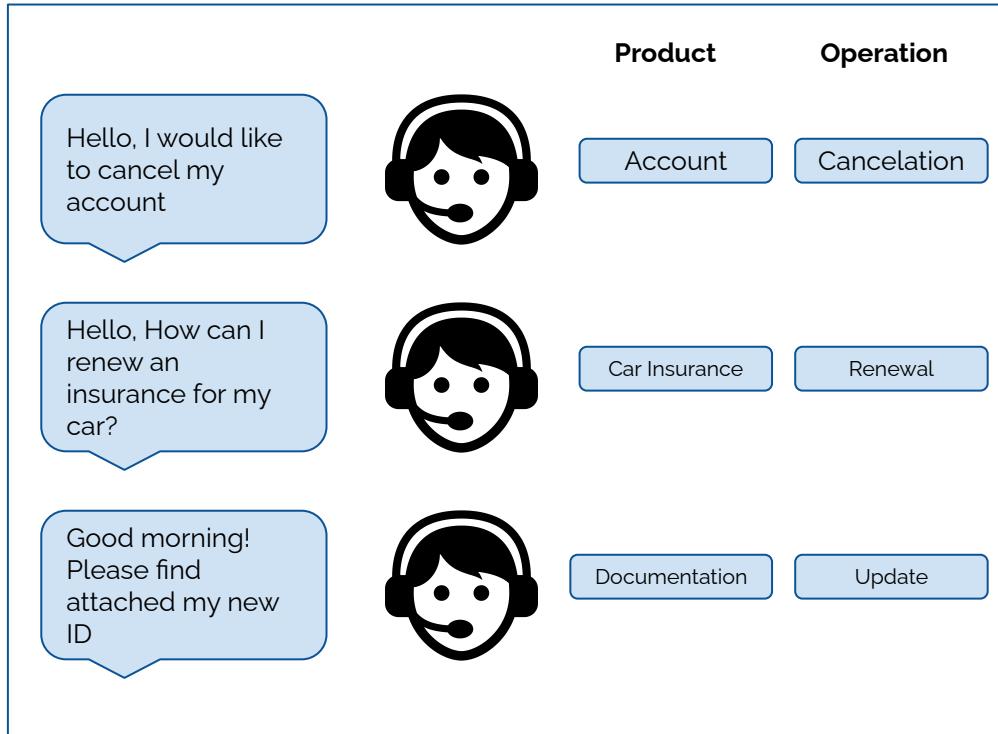
Post-modeling phases

- Deployment
- Monitoring
- Model updating



Real use case

Conversations classifier

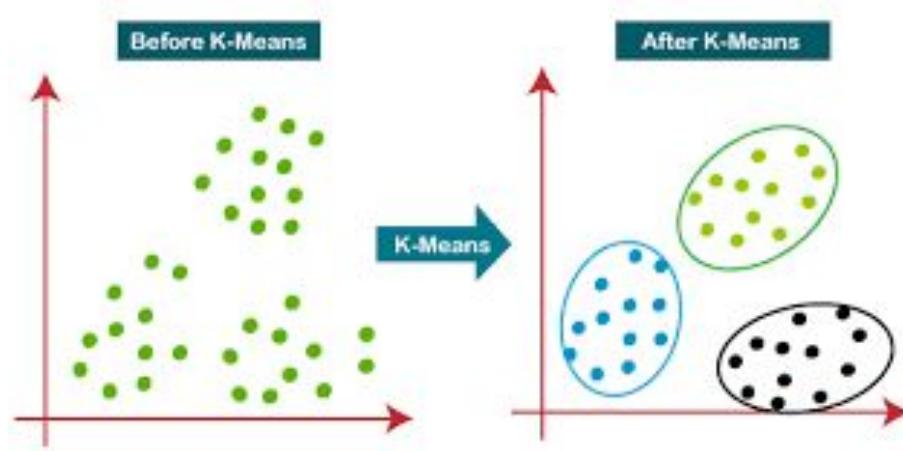


Unsupervised learning in NLP

Goal: Find structure inherent to unlabeled data

Common techniques in NLP

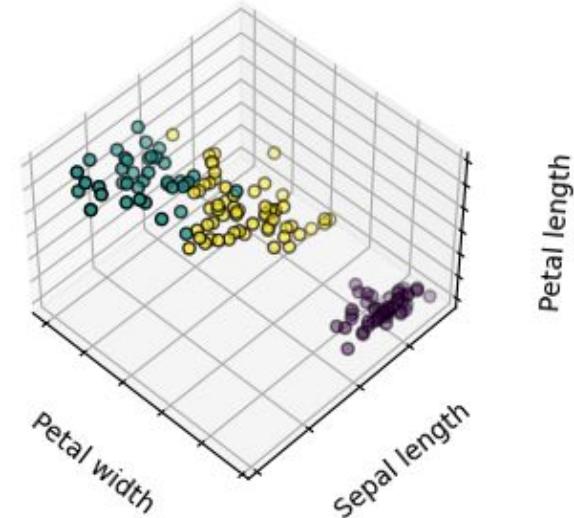
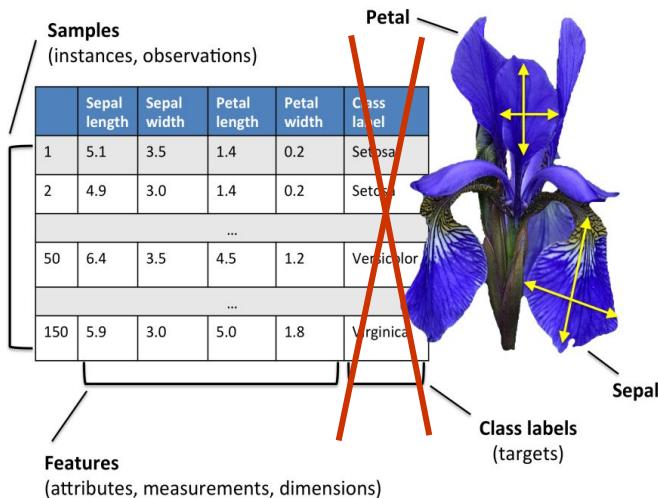
- Clustering
- Topic Modeling



Unsupervised learning in NL

Common techniques in NLP

- Clustering
- Topic Modeling



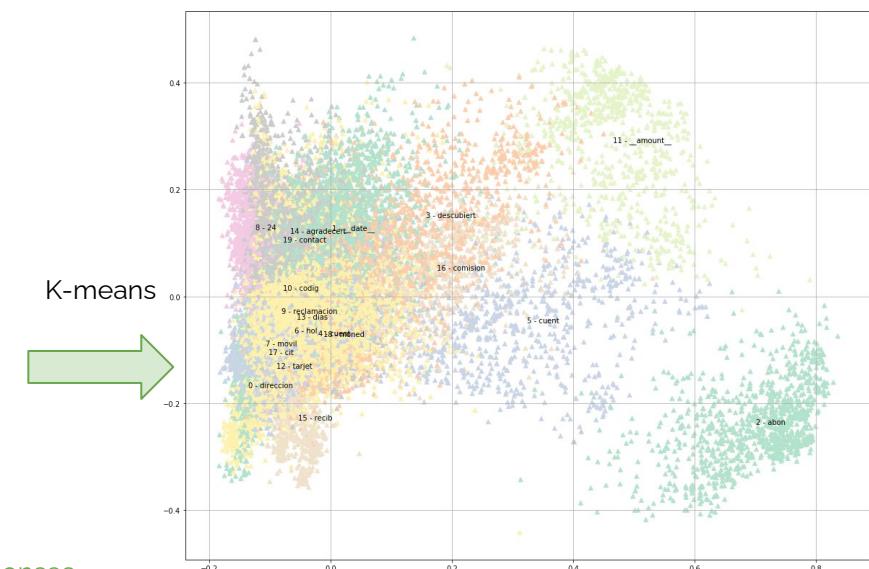
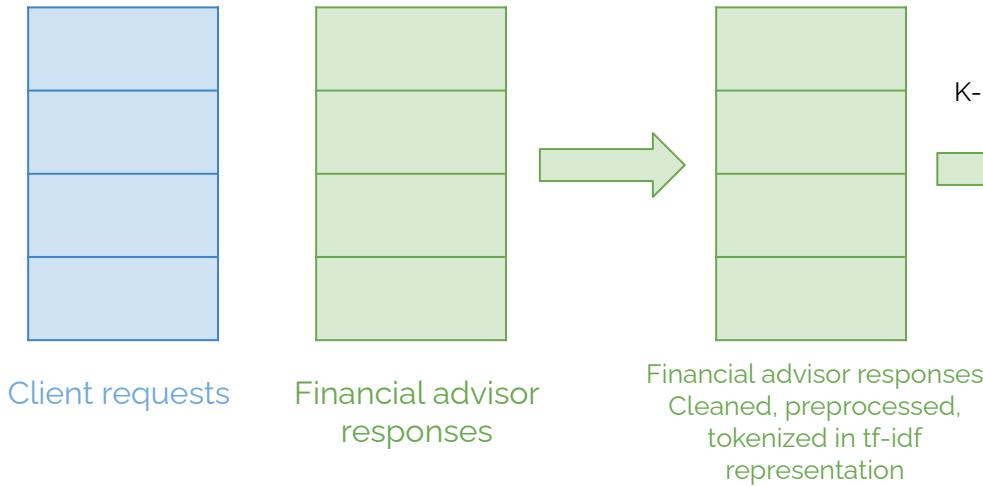
K-means

1. Input dataset in numerical representation
2. Select number of clusters
3. Initialize randomly centroid of each cluster
4. For each datapoint calculate distance to centroids and select minimum
5. Update centroid with new datapoint
6. Repeat step 4 until no data points are left

Unsupervised learning in NLP

Common techniques in NLP

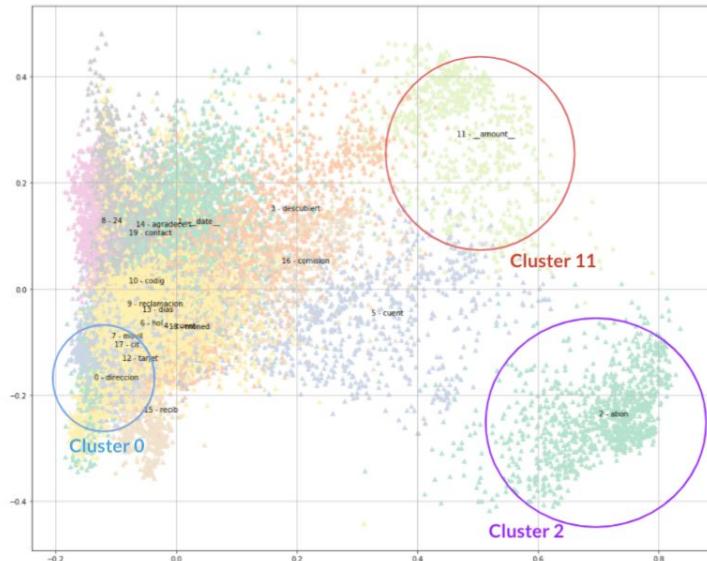
- Clustering



Unsupervised learning in NLP

Common techniques in NLP

- Clustering



Cluster 2: "Explicación de comisiones"

...se te han aplicado estos gastos al no cumplir los requisitos para estar dentro del PROGRAMA ADIÓS COMISIONES, en tu caso en particular es debido a que no tienes ingresos en **cuenta**. Para poder estar exento de **comisiones** de administración y mantenimiento, sería necesario cumplir con estos requisitos: **Abono** de nómina igual o superior a 600€ (mínimo 2 meses al **cuatrimestre**), o **Abono** por Desempleo o Pensión igual o superior a 300€ (mínimo 2 meses al **cuatrimestre**)....

...Las **comisiones**, en tu caso son: - Administración: **AMOUNT** (**AMOUNT** por apunte contable en **cuenta**). - Mantenimiento: **AMOUNT** (**AMOUNT** anuales)... **NAME**, actualmente no cumples con el programa "Adiós **comisiones**" en esta cuenta, pues, además de cumplir los requisitos de gastos ...

Cluster 11: "Comisiones por descubierto"

...Comisión de **descubierto** tácito: **HOUR_QR_AMOUNT** euros (3.6719% sobre el **importe** **NAME** de **descubierto** en el periodo de liquidación) por el **descubierto** por importe de **AMOUNT** producido el dia **DATE**...

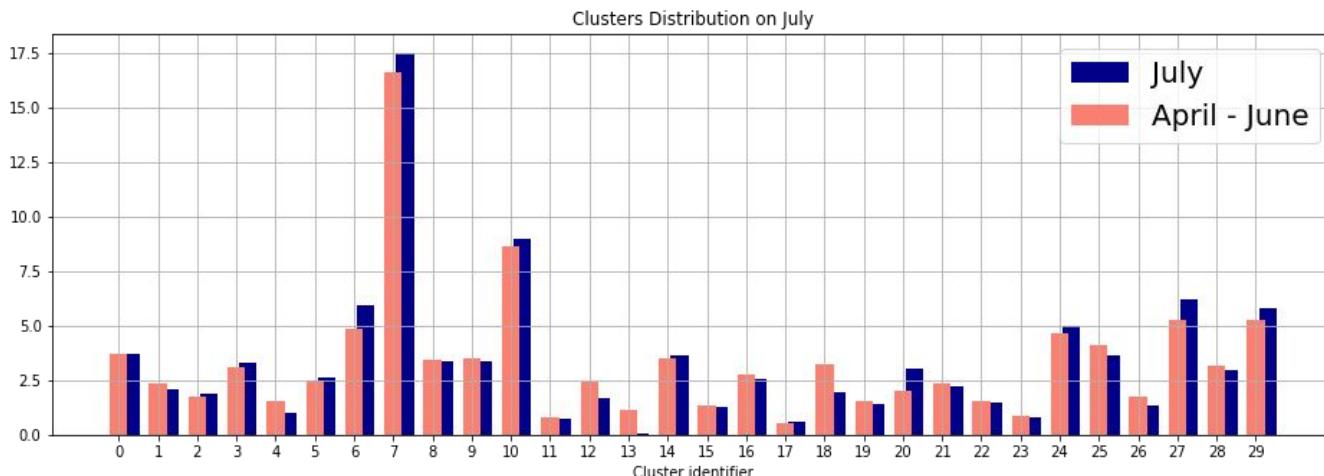
Cluster 0: "Cambio de domicilio fiscal"

... indicarte que puedes **modificar** la dirección postal de tus **productos** a través de la web ó la aplicación móvil de BBVA. A continuación te voy a indicar ambas rutas. **NAME**, para **modificar** el **domicilio** postal de tus **productos**...

Unsupervised learning in NLP

Common techniques in NLP

- Clustering



- Clustering model trained with April to June data
- July conversations clustered with the trained model

Unsupervised learning in NLP

Common techniques in NLP

- Clustering

Smart Reply: Automated Response Suggestion for Email

Anjuli Kannan* Karol Kurach* Sujith Ravi*

Tobias Kaufmann* Andrew Tomkins Balint Miklos

Greg Corrado László Lukács Marina Ganea

Peter Young Vivek Ramavajjalai

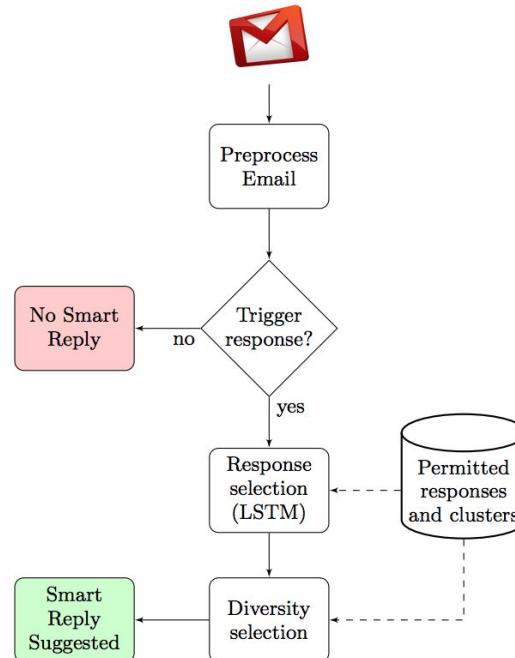
Google
{anjuli,kkurach,sravi,snufkin}@google.com

ABSTRACT

In this paper we propose and investigate a novel end-to-end method for automatically generating short email responses, called Smart Reply. It generates semantically diverse suggestions that can be used as complete email responses with just one tap on mobile. The system is currently used in *Inbox by Gmail* and is responsible for assisting with 10% of all mobile responses. It is designed to work at very high throughput and process hundreds of millions of messages daily. The system exploits state-of-the-art, large-scale deep learning.

We describe the architecture of the system as well as the challenges that we faced while building it, like response diversity and scalability. We also introduce a new method for semantic clustering of user-generated content that requires only a modest amount of explicitly labeled data.



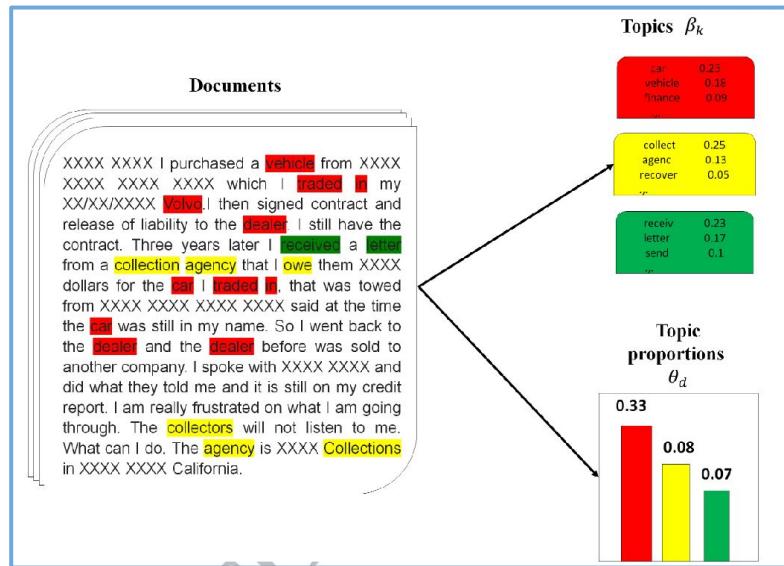


[Smart Reply: Automated Response Suggestion for Email](#)

Unsupervised learning in NLP

Common techniques in NLP

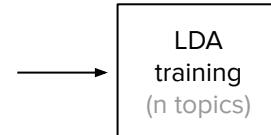
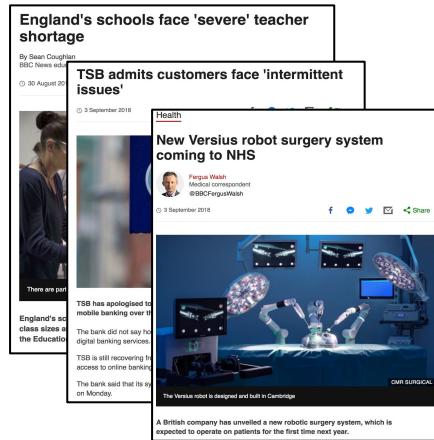
- Topic modeling
 - Use of a statistical model to discover the abstract "topics" that occur in a collection of documents.
 - Algorithms: **Latent Dirichlet Allocation**, Latent Semantic Allocation



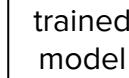
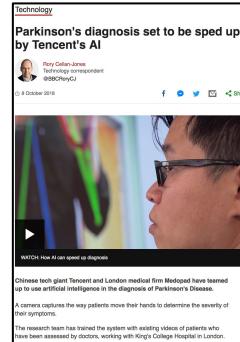
Unsupervised learning in NLP

Topic Modeling Latent Dirichlet Allocation (LDA)

- Generative probabilistic model
- Describes documents as collection of topics
- Applications
 - Document classification
 - Topic extraction



- Unsupervised: number of topics needed a priori
- Language agnostic
- Python using Gensim



Technology,
Health..

Topic Modeling Latent Dirichlet Allocation (LDA)

Document 0	politics parliament politics parliament
Document 1	cell biology cell biology
Document 2	biology cell
Document 3	parliament politics
Document 4	politics biology

	biology	cell	parliament	politics
Document 0	0	0	2	2
Document 1	2	2	0	0
Document 2	1	1	0	0
Document 3	0	0	1	1
Document 4	1	0	0	1
	biology	cell	parliament	politics

	Topic 0	Topic 1
biology	0.570	0.001
cell	0.428	0.001
parliament	0.001	0.428
politics	0.001	0.570
	Topic 0	Topic 1

	Topic 0	Topic 1
Document 0	0.071	0.929
Document 1	0.929	0.071
Document 2	0.875	0.125
Document 3	0.125	0.875
Document 4	0.500	0.500
	Topic 0	Topic 1

Topic Modeling Latent Dirichlet Allocation (LDA)

KL = 6.74

Applications: content similarity recommenders

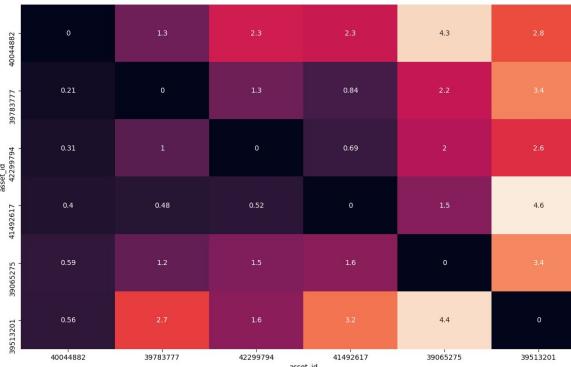


Dataset:
70,000 BBC articles published in 2017
English language

Docs	1	2	3	4	5	6	...
The Irish border Brexit backstop	0.7	0	0.2	0	0.1	0	
Scotland to get AI health research centre	0	0	0.9	0	0	0.1	

Similarity measurement

- Discrete probability distributions
- Kullback-Leibler divergence or relative entropy
 - Information gain between distributions



Topic Modeling Latent Dirichlet Allocation (LDA)

Applications: reading user profiles for segmentation

	women		International	Trump	Miscellaneous	Class	Companies	Road	Hospital	Murder/arrests	...
International	militant	gender	relations	travel	- celtic text	A	news	closures	problems	Murder/arrests	...
groups	groups	pay gap	upheaval	ban		drugs		and incidents			
0	attack	woman	president	university	belfast	drug	company	road	health	police	...
1	group	female	country	student	county	yorkshire	business	prison	hospital	man	...
2	kill	man	mr	ban	inquest	sheffield	firm	close	care	yearold	...
3	people	male	protest	law	ysgol	cannabis	pay	traffic	patient	arrest	...
4	military	gender	government	bill	qatar	supply	sell	m	nhs	attack	...
5	state	women	china	state	coroner	norfolk	m	crash	trust	woman	...
6	force	weinstein	state	new	primary	heroin	customer	bike	service	incident	...
7	paper	pay	iran	country	school	people	buy	reopen	treatment	find	...
8	muslim	equality	saudi	right	londonderry	find	sale	prisoner	staff	officer	...
9	violence	lee	chinese	immigration	st	suffolk	store	route	doctor	murder	...
...	report	work	israel	allow	paper	substance	money	near	cancer	hospital	...
...	army	harvey	minister	order	armagh	cocaine	share	cyclist	mental	street	...
...	soldier	singapore	opposition	people	cornwall	class	shop	driver	flight	road	...

Manual topic naming with BBC editors

Topic Modeling Latent Dirichlet Allocation (LDA)

Applications: reading user profiles for segmentation



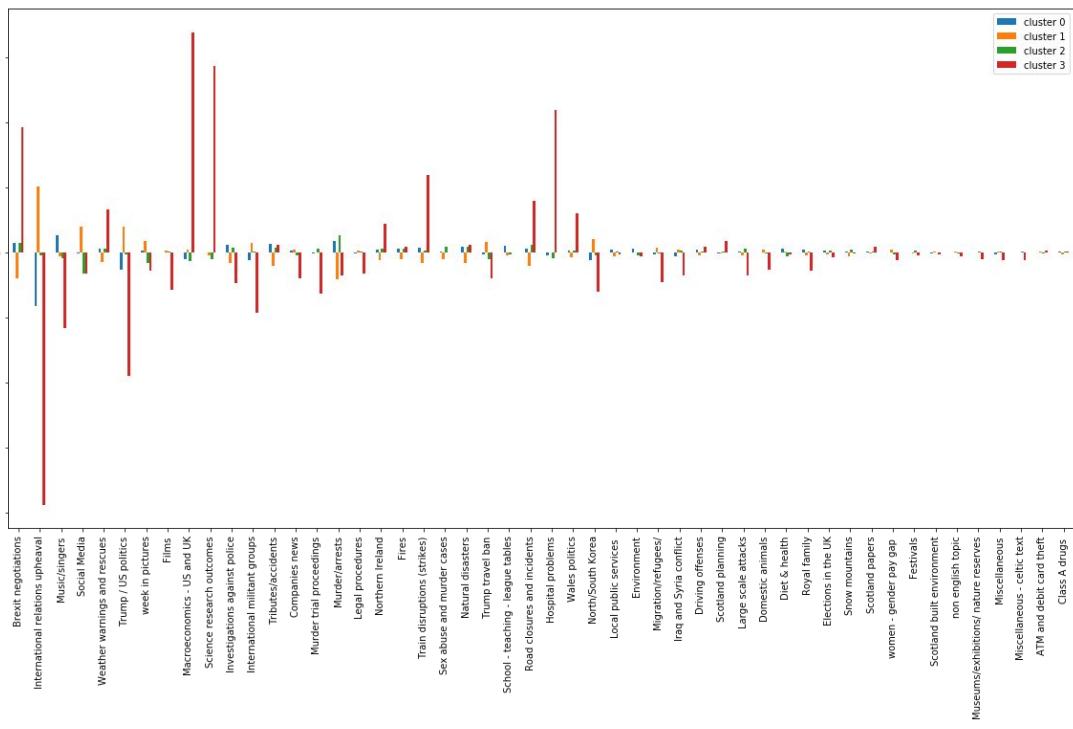
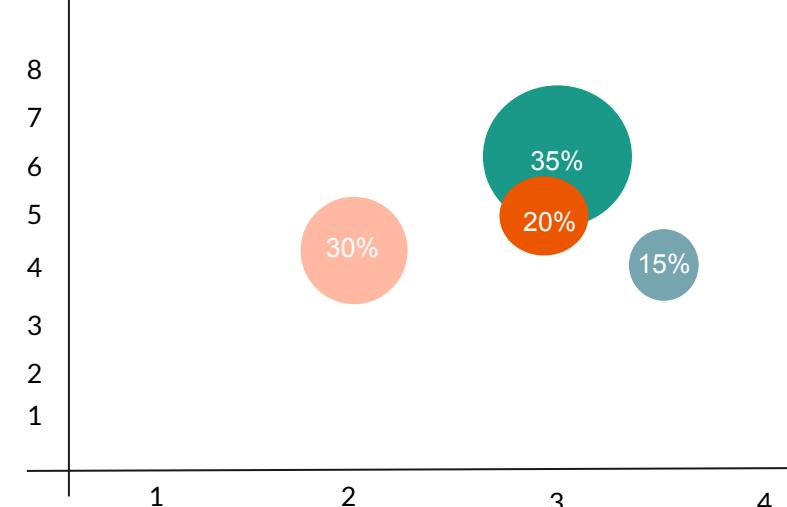
health	Scottish politics	education	local politics	brexit	crime	...
0.8	0	0.3	0	0	0.7	



health	Scottish politics	education	local politics	brexit	crime	...
0	0.6	0	0.5	0.6	0	

Topic Modeling Latent Dirichlet Allocation (LDA)

Applications: reading user profiles for segmentation



Topic Modeling Latent Dirichlet Allocation (LDA)

Applications: Political knowledge gaps

KNOWLEDGE GAPS AND THE ROLE OF THE BBC IN THE ERA OF FAKE NEWS

•••

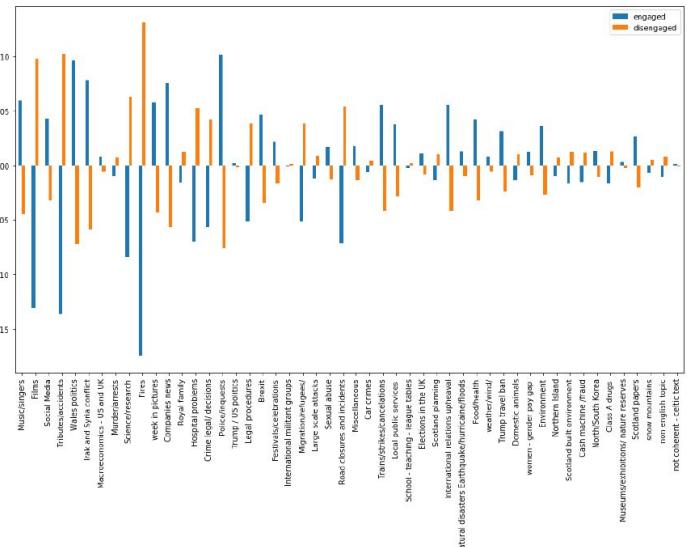
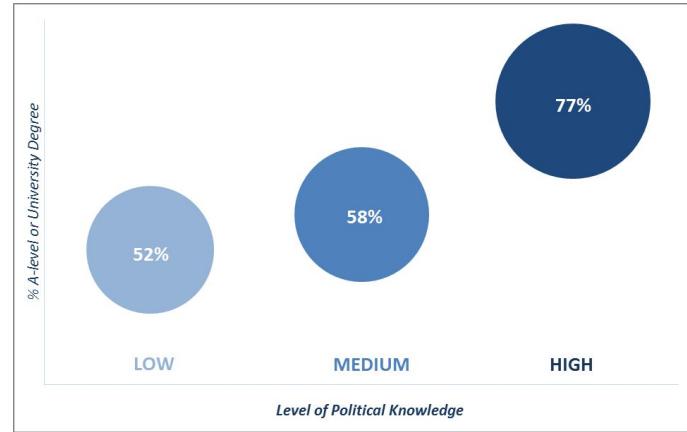
Marta Cantijoch

Nikki Soo

Tim Cowlishaw

Clara Higuera Cabañas

BBC R&D IRFS / BBC News Audience Engagement
University of Manchester Department of Politics



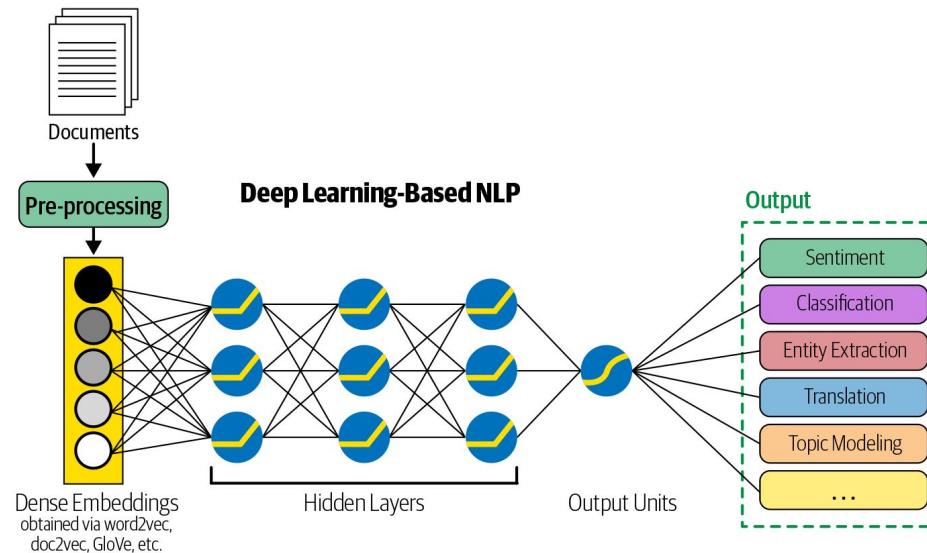
Hands on exercise...

NLP with deep learning

Preprocessing

Text representation

- Embeddings



Word embeddings

Word embeddings are a **type of word representation** that allows words that **share similar context** to have a similar representation.

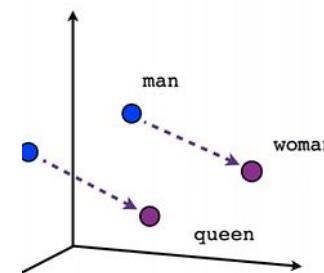
Use neural network architectures to create **dense, low-dimensional representations** of words and texts

Based on:

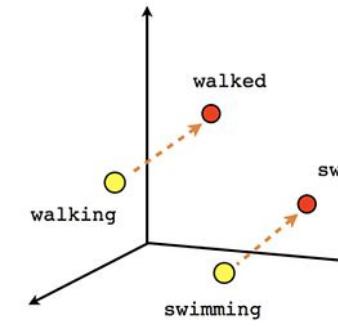
Distributional hypothesis

In linguistics, this hypothesizes that words that occur in similar contexts have similar meanings. For example, the English words "dog" and "cat" occur in similar contexts.

Thus, if two words often occur in similar context, then their corresponding representation vectors must also be close to each other.



Male-Female

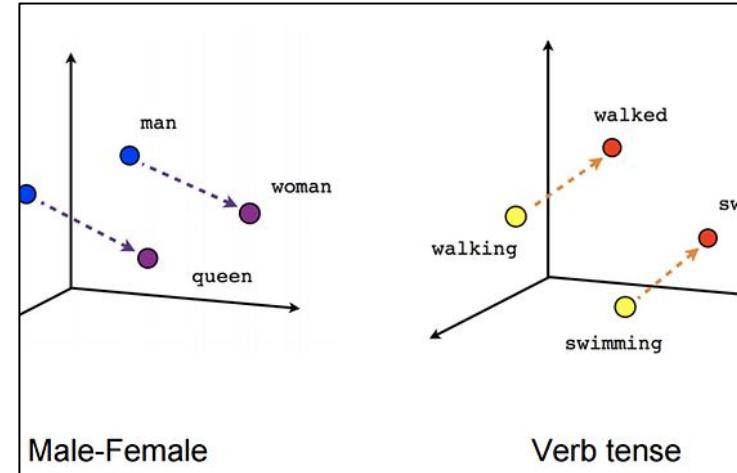


Verb tense

Word embeddings

Neural network-based word representation model known as "**Word2vec**," based on "distributional similarity," can capture word analogy relationships such as

$$\text{King} - \text{Man} + \text{Woman} \approx \text{Queen}$$



- Captures context in vector representation vs BOW /tf-idf
- Fixed low dimensionality vectors vs BOW/tf-idf increasing vectors with increased vocabulary

Word embeddings

1. Select a embedding dimension and a window size (hyperparameters)
2. For every word w in corpus, we start with a vector v_w initialized with random values.
 - a. **CBOW model:** predicts the center word given the context words in which the center word appears.
 - b. **Continuous Skipgram model:** predict the context words from the center word.

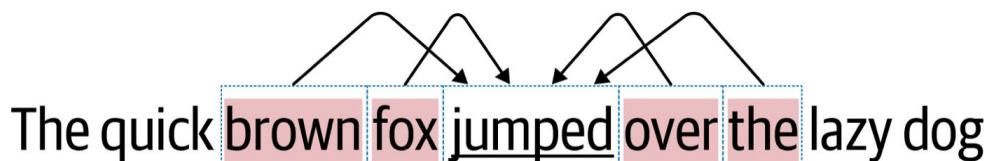


Figure 3-7. CBOW: given the context words, predict the center word

Source Text

Training Samples
(context, target)

The quick brown fox jumps over the lazy dog.	→	((quick, brown), The)
The quick brown fox jumps over the lazy dog.	→	((The, brown, fox), quick)
The quick brown fox jumps over the lazy dog.	→	((The, quick, fox, jumps), brown)
The quick brown fox jumps over the lazy dog.	→	((quick, brown, jumps, over), fox)

Figure 3-8. Preparing a dataset for CBOW

Word embeddings

CBOW Model

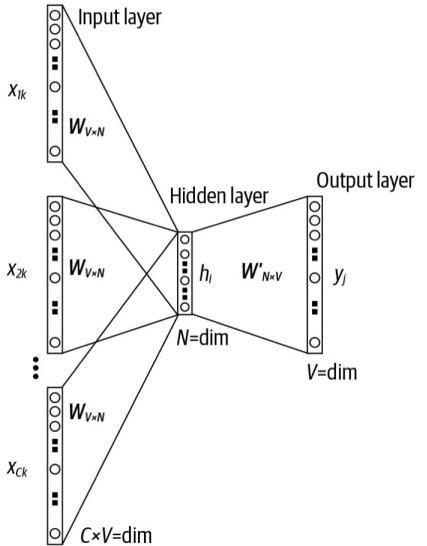
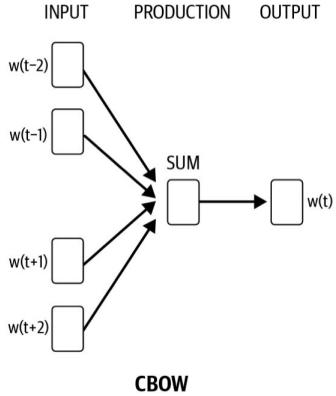


Figure 3-9. CBOW model [7]

[6] #What is the vector representation for a word?
w2v_model['beautiful']

```

array([-0.03831055,  0.0566466, -0.01153564,  0.07275391,  0.15136719,
       -0.0495152,   0.24851562, -0.09472656, -0.03344727,  0.24797031,
       -0.0644512,   0.24851562, -0.09472656, -0.03344727,  0.24797031,
       0.05541931, -0.08921631,  0.1328125, -0.15429688,  0.08185469,
      -0.07373047,  0.24316466,  0.12353516, -0.09277344,  0.08283125,
      0.06494141,  0.15722656,  0.11279297, -0.0612793, -0.296875,
      -0.11279297,  0.24316466,  0.12353516, -0.09277344,  0.08283125,
      -0.27539062,  0.0644512,   0.13867188, -0.08897318,  0.09472656,
      0.07861328, -0.10058594,  0.23925781,  0.03888594,  0.18652344,
     -0.11279297,  0.22558594,  0.10986328, -0.11865234,  0.02026367,
     0.11376953,  0.09579312,  0.29492188,  0.08251953, -0.05444336,
      -0.04746332, -0.08544922,  0.10942969, -0.30273438,  0.07617189,
      0.125, -0.05932617,  0.03833808, -0.03564453,  0.2421875,
     0.36132812,  0.04760742,  0.00631714, -0.03988379, -0.13964844,
     0.22558594, -0.06298828, -0.02636719,  0.1171875,  0.33398438,
     -0.07656161,  0.06869453,  0.0093891, -0.16091369, -0.22449338,
     0.02026367,  0.0332266,  0.1712081, -0.16991108, -0.1572266,
     0.13865938,  0.12451172, -0.20418156,  0.04736328, -0.296875,
     -0.17488469,  0.00872803,  0.04638672,  0.10791816, -0.283125,
     -0.27539062,  0.2734375, -0.02563477, -0.11835156,  0.0625,
      0.1953162,  0.16846262, -0.1376953, -0.03636363, -0.095375,
     0.22131562,  0.0390625, -0.03894424, -0.03894424,  0.309375,
     0.1796875,  0.03833808, -0.24804688,  0.03515625,  0.03881836,
     0.03442383, -0.04101562,  0.20214844, -0.03015137, -0.09619141,
     0.11669922, -0.06738281,  0.0625,  0.10742188,  0.25585938,
     -0.21777344,  0.05639645,  0.00938918, -0.16113281,  0.11865334,
     0.00832018,  0.0332266,  0.1712081, -0.16991108, -0.1572266,
     0.05883789,  0.00634766,  0.11914062,  0.07324219, -0.01586914,
     0.18457931,  0.0532266,  0.19824219, -0.22363281, -0.25195312,
     0.15059962,  0.22753996,  0.05737305,  0.16992188, -0.22558594,
     0.06494141,  0.11914062, -0.06646625, -0.18449219, -0.07226562,
     -0.03833808,  0.03833808, -0.03894424, -0.03894424,  0.02172952,
     -0.1269512,  0.18457931,  0.27539962, -0.36523438,  0.03491211,
     -0.18554688,  0.23828125, -0.13867188,  0.08296021,  0.04272461,
     0.13867188,  0.12207031,  0.05957031, -0.22167969, -0.18945312,
     -0.23242188, -0.28719933,  0.08866699, -0.16113281, -0.24316406,
     0.00938918,  0.0332266,  0.09883464, -0.08866699, -0.08866699,
     0.16113281,  0.11621694,  0.16462025, -0.03697754,  0.01574707,
     0.11425781, -0.04174865,  0.05980285,  0.02661133, -0.08642578,
     0.140625,  0.09228516, -0.25195312, -0.31445312, -0.05688477,
     0.01031494,  0.0234375, -0.02331543, -0.08956641,  0.01269531,
     -0.1394688,  0.17742156, -0.03894424, -0.03894424,  0.02172952,
     0.11962808,  0.18457931, -0.20137178, -0.03140918,  0.02172952,
     0.23946875,  0.28125, -0.17675781,  0.02378516,  0.08748234,
     -0.06176758,  0.00939941, -0.09277344, -0.029125,  0.13885938,
     -0.13671875, -0.00509488,  0.04296875,  0.12988281,  0.3515625,
     0.00832018,  0.12207031, -0.03173828,  0.28515625,  0.18261119,
     0.13671875,  0.05308566,  0.09883464, -0.08866699, -0.08866699,
     0.08740234,  0.04021143, -0.1320125, -0.17578125, -0.04321289,
     -0.015625,  0.16894531,  0.23,  0.37109375,  0.19921875,
     -0.36132812, -0.10302734, -0.20800781, -0.20117188, -0.01519775,
     -0.12207031,  0.12817178, -0.07421875, -0.04345703,  0.14160156,
     -0.11376953,  0.08957031, -0.03894424, -0.03894424,  0.02172952,
     -0.03027344,  0.01796888, -0.10858594, -0.20793125,  0.11376953,
     -0.12402344,  0.04003986,  0.06933594, -0.34579312,  0.03881836,
     0.16210938,  0.045761719, -0.12792969, -0.05810547,  0.03857422,
     0.11328125, -0.1953125, -0.2285, -0.13183594,  0.15722656,
     0.00961914,  0.09966938, -0.00285339, -0.03637695,
     0.15429688,  0.06152344, -0.34579312,  0.11083904,  0.03347271,
  dtype='float32')

```

Figure 3-6. Vector representing the word “beautiful” in pre-trained Word2vec

Word embeddings

Ways of use

Pretrained models

- [Word2vec in gensim](#)
- [Glove - Stanford](#)
- [FastText\(Facebook\)](#)

Train your own model

- Useful when your use case is context specific

```
#Import a test data set provided in gensim to
train a model

from gensim.test.utils import common_texts
#Build the model, by selecting the parameters.
our_model = Word2Vec(common_texts, size=10,
window=5, min_count=1, workers=4)
```

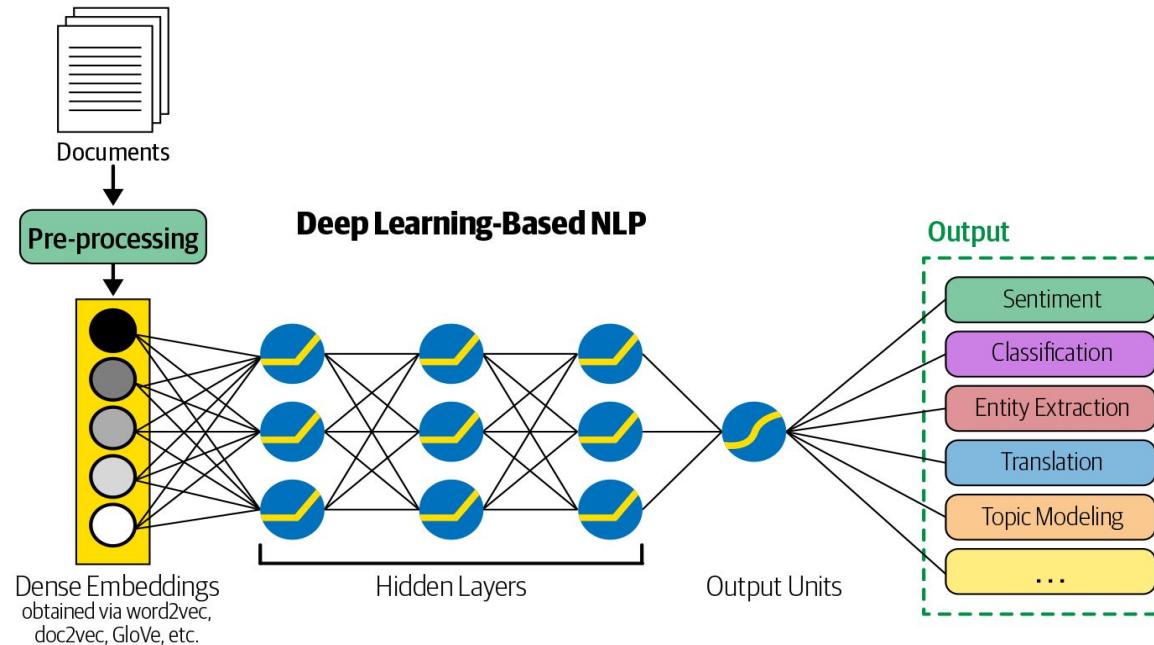
Chapter 3. Text Representation.[Practical Natural Language Processing. O'Reilly](#)

[Word2vec tutorial](#)

[CBOW model architecture - video](#)

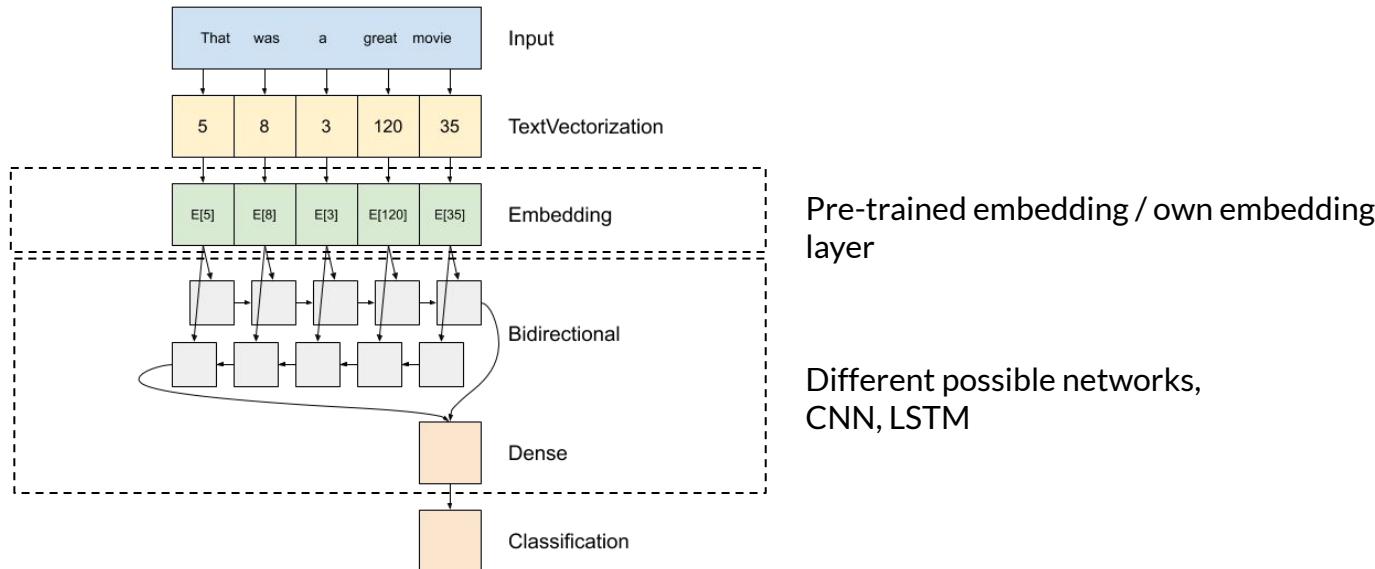
Hands on exercise...

NLP with deep learning



Text classification with deep learning

Text classification with deep learning

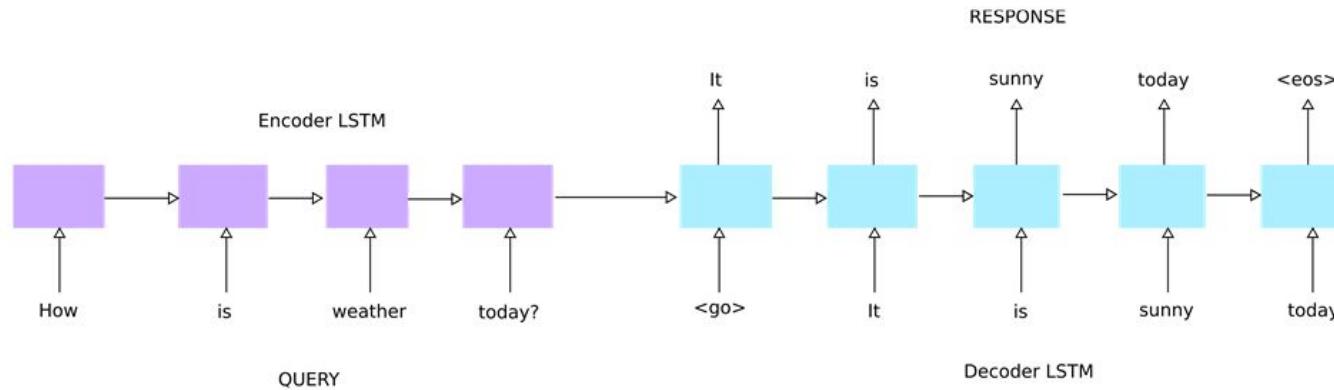


[Practical example notebook](#)

Natural language generation

Applications

- Machine translation
- Summarization
- Question Answering



Question Answering with seq2seq

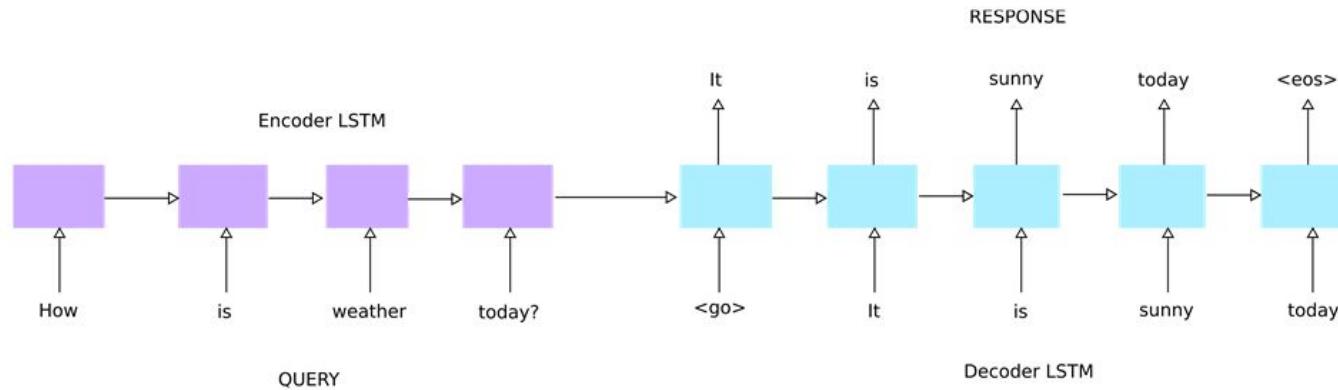


Natural language generation

Applications

- Machine translation
- Summarization
- Question Answering

[Video: Prof. John Kelleher explains neural machine translation](#)

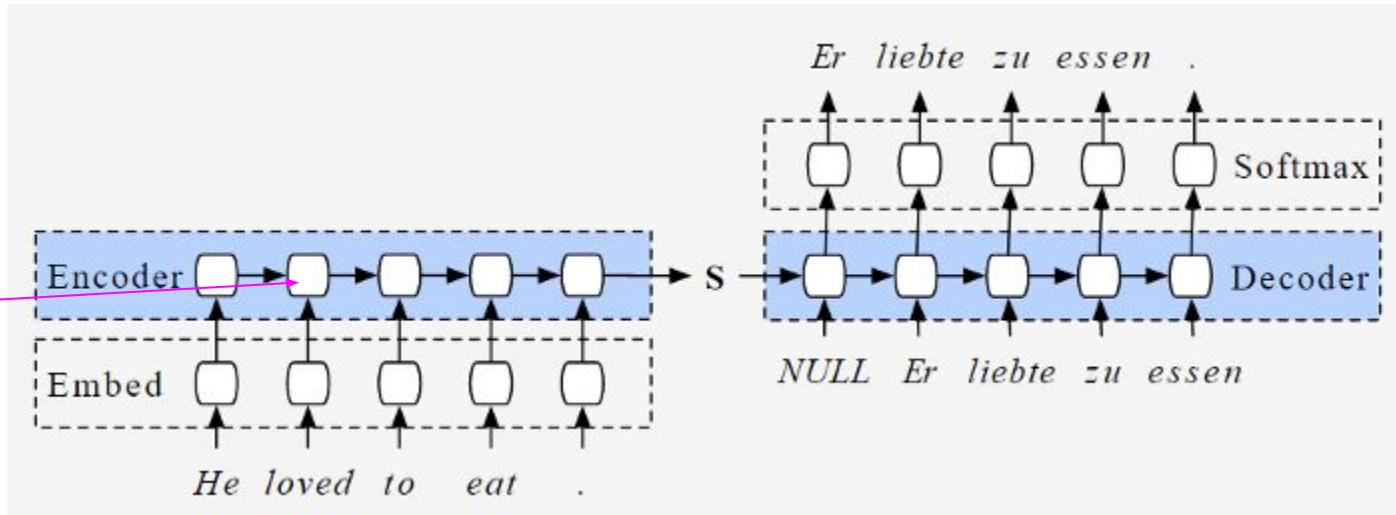


Question Answering with seq2seq

Natural language generation

seq2seq

hidden states: in
training capture info
about the current cell
and previous one



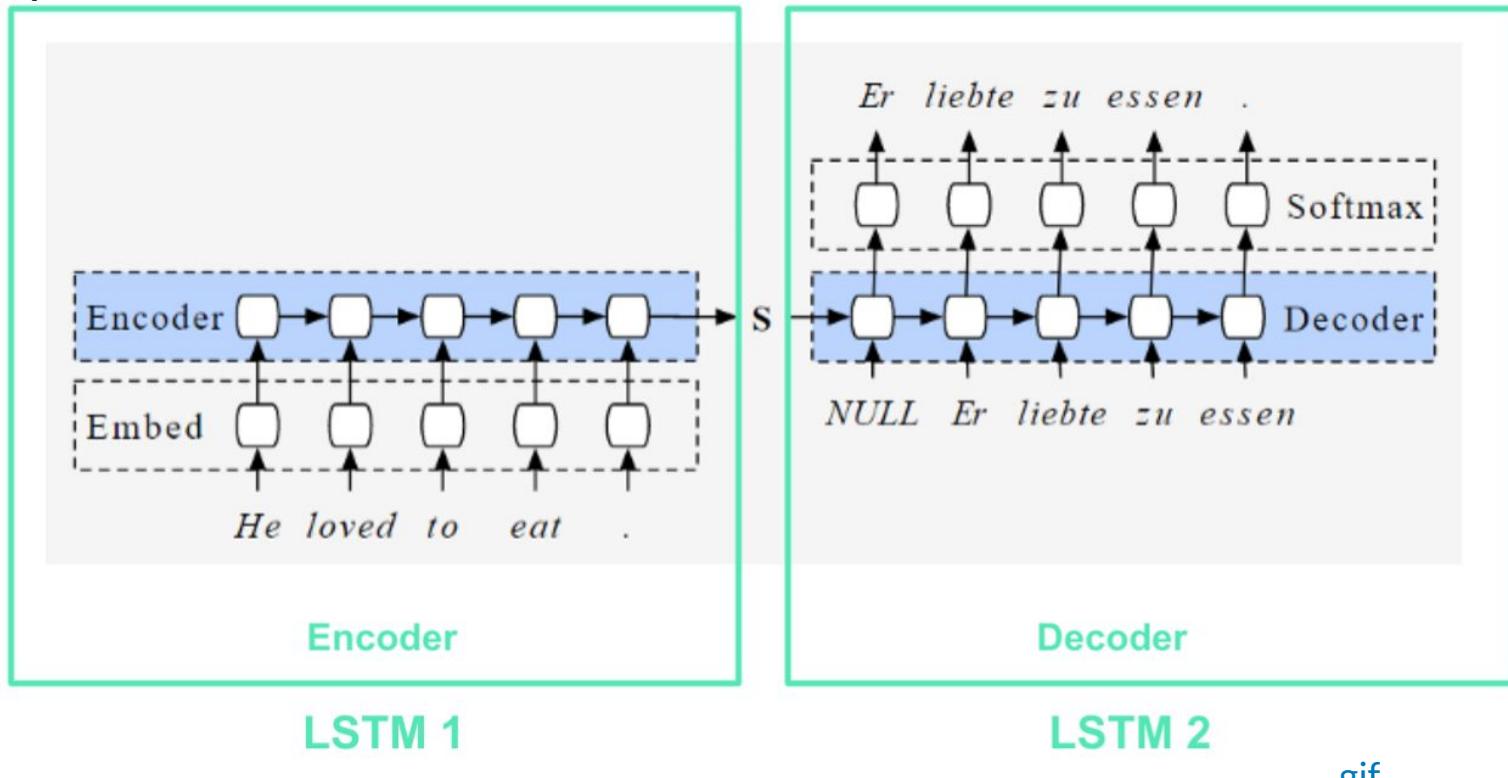
Encoder

Vector that
represents entry
sequence

Decoder

Natural language generation

seq2seq

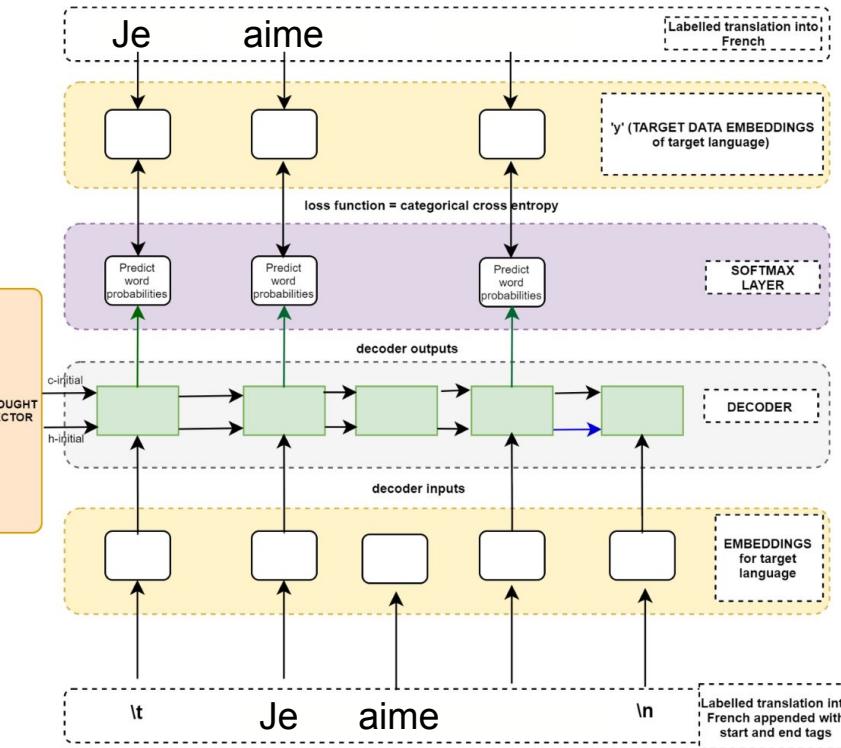
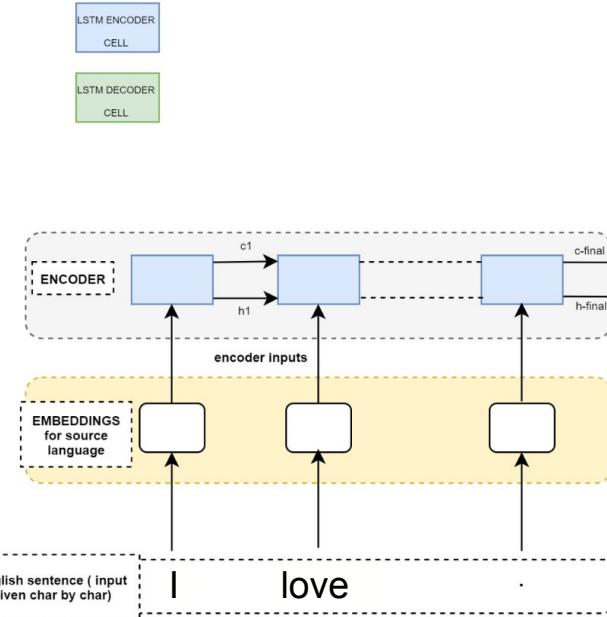


[gif](#)

Natural language generation

seq2seq

ENCODER - DECODER TRAINING NETWORK
ARCHITECTURE FOR NEURAL MACHINE
TRANSLATION



Training

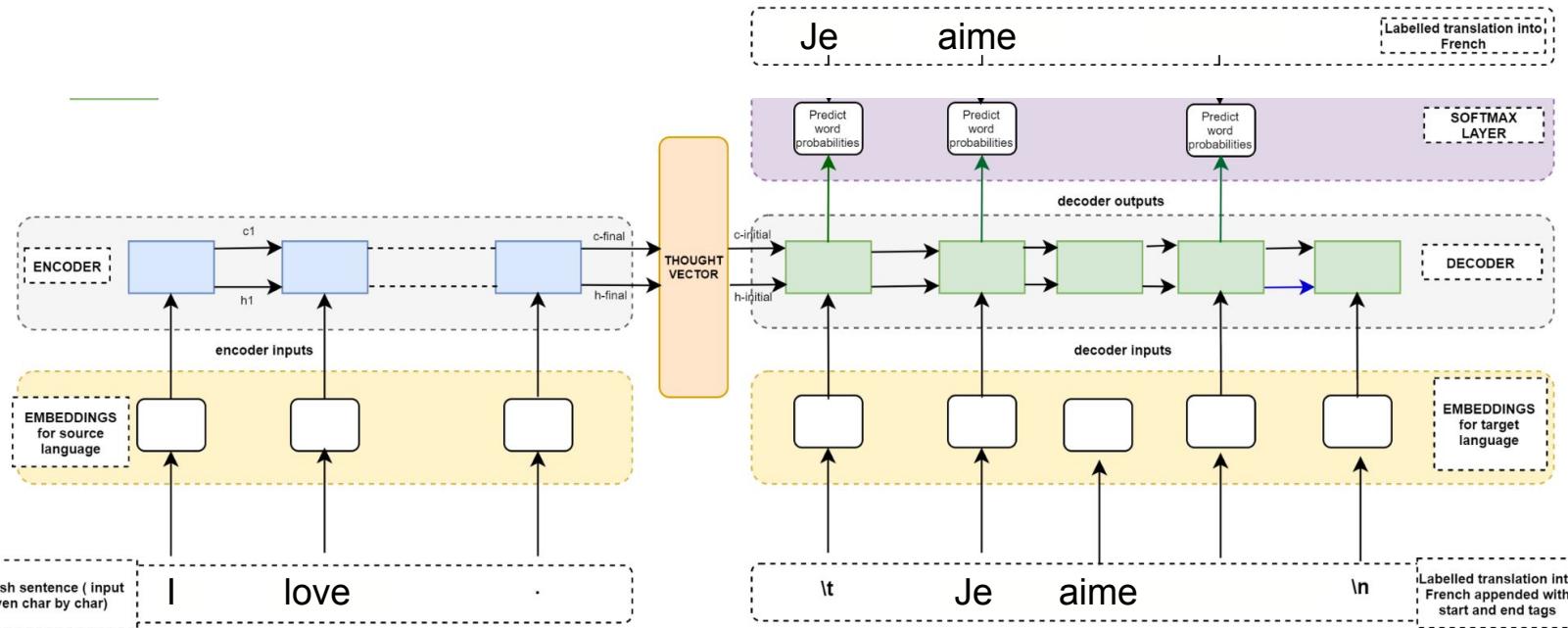
Parameters

- max_len_seq
- num epochs
- embedding_dim
- learning rate
- attention layer
- drop out
- ...

Natural language generation

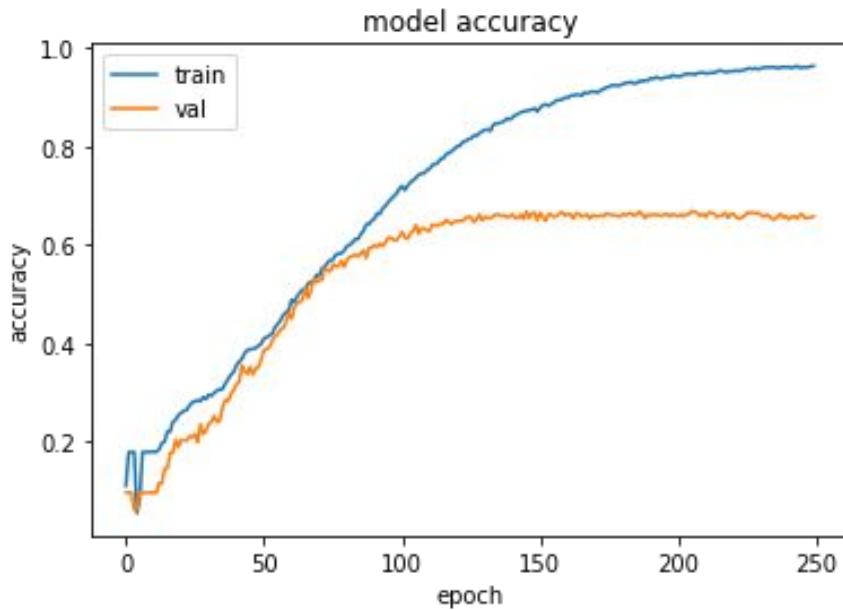
Inference

seq2seq



Natural language generation

NLG Evaluation



respuesta real

saludo	Para	solicitar	una	nueva	tarjeta	debes
saludo	Debes	solicitar	la	nueva	en	la

respuesta generada automáticamente

Natural language generation metrics
(bleu, rouge, f1...)

Natural language generation

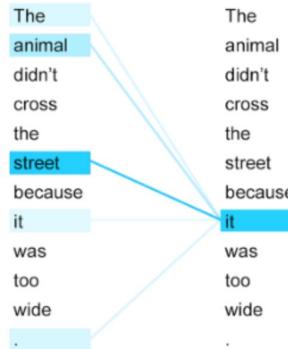
Limitations of seq2seq

- Dealing with long-range dependencies is still challenging
- The sequential nature of the model architecture prevents parallelization.

The **Transformer in NLP** is a novel architecture that aims to solve sequence-to-sequence tasks while handling long-range dependencies with ease.

[Attention Is All You Need.](#)

“Self-attention, sometimes called intra-attention, is an attention mechanism relating different positions of a single sequence in order to compute a representation of the sequence.”



[Understanding transformers](#)

What's left...

1. Transfer learning in NLP
2. Explainability in NLP
3. Bias and dangers in NLP

What's left...

1. Transfer learning in NLP
2. Explainability in NLP
3. Bias and dangers in NLP

y=python (probability 0.005, score -7.792) top features

Contribution?	Feature
-3.737	<BIAS>
-4.054	Highlighted in text (sum)

```
add element from a list of dictionary to another <pre><code>dict_1 = [ { incident_id : sd000001372596 first_call : t } { incident_id : sd000001372594 first_call : f } { incident_id : sd000001372598 first_call : f } { incident_id : sd000001372599 first_call : t } { incident_id : sd000001372602 first_call : f } { incident_id : sd000001372601 first_call : f } { incident_id : sd000001372605 first_call : f } { incident_id : sd000001372606 first_call : f } { incident_id : sd000001372607 first_call : f } ] dict_2 = [ { incident_id : sd000001372605 date : 08-10-2016 00:54:13 } { incident_id : sd000001372606 date : 08-10-2016 00:57:20 } { incident_id : sd000001372607 date : 08-10-2016 01:00:25 } { incident_id : sd000001372598 date : 11-10-2016 10:57:34 } { incident_id : sd000001372602 date : 08-10-2016 10:44:34 } { incident_id : sd000001372601 date : 21-10-2016 22:30:49 } { incident_id : sd000001372594 date : 18-10-2016 14:53:34 } ] </code></pre> i have two list of dictionaries with different length and i want to add the ( dict_2 date ) to the dict_1 according the incident_id
```

Eli5 library

y=jquery top features		y=mysql top features		y=objective-c top features		y=php top features	
Weight?	Feature	Weight?	Feature	Weight?	Feature	Weight?	Feature
+3.977	jquery	+4.897	mysql	+4.001	objective	+2.558	php
+2.060	ready	+1.316	million	+1.276	passcode	+1.607	ga
+1.166	jsfiddle	+1.158	isregistered	+1.268	monish	... 23131 more positive ...	
+1.109	fadein	+0.962	restore	+1.262	cocoa	... 128982 more negative ...	
+1.090	animate	+0.927	mdf	+1.241	zziplib	-1.271	iphone
... 24776 more positive ...		+0.920	truncated	... 35356 more positive ...		-1.278	java
... 127337 more negative ...		+0.905	tablename	... 116757 more negative ...		-1.282	objective
-1.029	blockquote	... 32999 more positive ...		-1.155	who	-1.283	page_title_home
-1.035	employee	... 119114 more negative ...		-1.182	comment	-1.339	void
-1.056	scope	-0.921	connections	-1.204	def	-1.340	flags
-1.081	python	-0.939	python	-1.206	stringappendformat	-1.366	logout
-1.087	folder	-0.954	app	-1.223	jquery	-1.379	fields_by_id
-1.087	int	-0.992	web	-1.232	mysql	-1.398	javascript
-1.151	mean	-1.001	product_table	-1.243	selectedstring	-1.450	clicking
-1.164	sql	-1.033	2010	-1.282	textfielddidendediting	-1.566	rails
-1.192	javascript	-1.065	define	-1.422	rails	-1.600	net
-1.333	getelementbyid	-1.113	conn	-1.437	css	-1.617	angularjs
-1.502	rails	-1.124	net	-1.471	net	-1.705	footer
-1.650	angularjs	-1.161	name_table	-1.497	javascript	-1.711	uri
-1.930	ng	-1.642	android	-1.631	android	-2.092	python
-2.222	angular	-1.889	rails	-1.749	swift	-2.463	jquery
-3.738	<BIAS>	-3.892	<BIAS>	-2.530	<BIAS>	-3.317	<BIAS>

What's left...

1. Transfer learning in NLP
2. Explainability in NLP
3. **Bias and dangers in NLP**

Bias in Natural Language Processing (NLP): A Dangerous But Fixable Problem

One of the biggest new issues that natural language processing (NLP) models face is the implicit biases that they learn.

<https://towardsdatascience.com/bias-in-natural-language-processing-nlp-a-dangerous-but-fixable-problem-7d01a12cf0f7>

References

Books

[Practical Natural Language Processing, O'Reilly](#)

[Practical Natural Language Processing with Python with Case Studies from Industries Using Text Data at Scale](#)

[Generative Deep Learning. Teaching Machines to Paint, write Compose and Play. O'Reilly](#) (Chapter 6 Language generation)

Interesting articles/posts

[How AI is getting better at Detecting Hate Speech](#)

[The dark side of Guardian Comments](#)

Scientific papers

[Smart Reply: Automated Response Suggestion for Email](#)