

BACK TO BASICS NLP

Octubre 2018
PyConEs Málaga
Claudia Guirao

- NATURAL LANGUAGE PROCESSING -

notes for PyConES 2018
@claudiaquirao

①

TEXT CLEANING

- regular expressions (regexp): pattern, looking for regular text, useful for identifying removing extracting text entities

- text/string transformations with



- what to do with punctuation?
 - conserve → emotions !!
 - remove → standardize
 - what about emojis? 😱

- dealing with the "encoding nightmare" UTF-8 ASCII

②

TEXT STANDARDIZATION

- tokenization: split our text into tokens a.k.a. words

- stemming: take the stems of our tokens

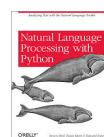
some stemmers

none of them works so well on Spanish

- PORTER
- SNOWBALL
- LOVINS
- PALICE

How to perform this task?

custom functions



spaCy
FreeLing

CLiPS
COMPUTATIONAL LINGUISTICS & PSYCHOLINGUISTICS RESEARCH CENTER

- lemmatization: grouping together the inflected forms so they can be analysed as a single item

- make decisions about other expressions



spatially relevant on informal contexts

LOL
WTF
IMHO

- STOPWORDS: most common words in a language.

↑ a universal stop words list → it depends of the domain

TIP build your own list and iterate

some libraries has its own list of stop words

- NLTK
- SPACY
- Google STOPWORDS

III

MAIN NLP ELEMENTS

notes for PyConES 2018
@claudia.guirao

Corpus: large collection of documents

<ul style="list-style-type: none"> · monolingual · multilingual 	<ul style="list-style-type: none"> · general domain (news, books, etc.) · specific domain → search for corpus on your domain
---	--

· some libraries provide corpus on several fields → NLTK

I need a Spanish corpus on my analysis

general
specific

I need a Spanish tagged corpus on my analysis → good luck

↳ you will probably
end tagging manually your documents

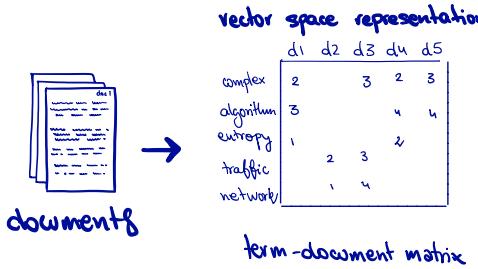


Documents → each unit in our corpus

- phrase?
- book?
- dialog?

each doc is made up of words = terms

IV BAG OF WORDS



Term frequency: measures how frequently a term occurs in a document

$$TF(t) = \frac{\# \text{ times term } t \text{ appears in a doc}}{\# \text{ of terms in the doc}}$$

Variants of term frequency (tf) weight	
weighting scheme	tf weight
binary	0, 1
raw count	$f_{t,d}$
term frequency	$f_{t,d} / \sum_{t' \in d} f_{t',d}$
log normalization	$\log(1 + f_{t,d})$
double normalization 0.5	$0.5 + 0.5 \cdot \frac{f_{t,d}}{\max_{\{t' \in d\}} f_{t',d}}$
double normalization K	$K + (1 - K) \frac{f_{t,d}}{\max_{\{t' \in d\}} f_{t',d}}$

Inverse document frequency: how important a term is.

TF considers all terms equally important

IDF weigh down the frequent terms while scale up the rare ones

$$IDF = \log_e \left(\frac{\text{total number of documents}}{\#\text{ of documents with term } t \text{ in it}} \right)$$

TF IDF = term "specificity"

combination of the freq. in the document and its freq. on the corpus

$$TFIDF(t, d, D) = TF(t, d) \cdot IDF(t, D)$$

how to perform these tasks?

1) custom counter functions

2) sklearn vectorizer.



documents → observations
sentiment / classification → target
terms → features

DISTANCES = useful for speech correction, search engines and analysis. We can calculate distances through 2 words by several distance measures

LEVENSHTEIN = distance between 2 words

WAGNER & FISHER MINIMUM EDIT DISTANCE = n Levenshtein distance
 insertion 1
 deletion 1
 replacement 2

INTENTION ↓
 EXECUTION
 1) delete i
 2) replace n by e
 3) substitute t by x
 4) insert u
 5) replace n by c

V DISTRIBUTIONAL REPRESENTATION

WORD₂VEC

- it is used to predict the nearby word
- in this case the order is relevant → the context
- based on NN
- 2 outputs
 - model for the next word
 - vectors that can be subtracted / added → represents the meaning of a word
- 2 ways to train this NN (CBOW // skipgram)

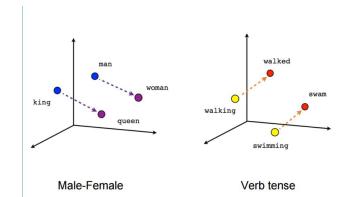
word2vec
 "PS! Thank you for such an awesome top"

the meaning of a word is given by its context

gensim
 topic modelling for humans
 unsupervised algorithm

the cat climbed a tree

Given context:
 a, cat, the, tree
 Estimate prob. of:
 climbed



DOC_1732 LDA

"PS! Thank you for such an awesome top"

LDA : LATENT DIRICHLET ASSIGNMENT unsupervised algorithm
 topic modelling algorithm, by Andrew NG & others

- works globally, it is a probabilistic model, learns a document vector that predicts words inside that document
- outputs a previously fixed number of topics and the words related with among their weights.

	0	1	2
Document 0	10	0	0
Document 1	0	10	0
Document 2	0	0	10
Document 3	10	10	10

Document 0		
Document 1		
Document 2		
Document 3		

LDA₂VEC mix both approaches -local -global unsupervised algorithm

from Chris Moody paper

VII USE CASES

notes for PyConES 2018
@claudiaquirao

1) LABELED DOCUMENTS

- sentiment analysis
- document classification

2) UNLABELED DOCUMENTS

- segmentation / clustering
- topic modelling
- entities recognition

other analysis & information retrieval

- indexers & search engines
- speech correction
- comparing with general domain
- ontology and network analysis
- translation
- tagging
- intents identification and chatbots

