

91258 - Natural Language Processing

Lesson 3. Vector Space Model

Alberto Barrón-Cedeño

Alma Mater Studiorum-Università di Bologna
a.barron@unibo.it @albarron_

09/10/2023



Table of Contents

Current Status

Representations Revisited

More Basic Algebra

Sentiment Analysis

Current Status

Current Status

You know...

- ▶ what is natural language processing
- ▶ there are two main paradigms: rule-based and statistical

On your own, you have...

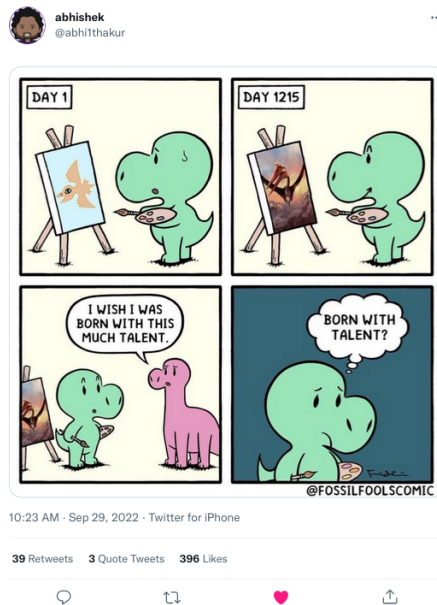
- ▶ setup a Python development environment
 1. command line
 2. PyCharm or any other option (e.g., Eclipse)
 3. Google's Colab

On your own, you (could) have...

- ▶ found out what is **git** (and perhaps \LaTeX as well!)

You can...

- ▶ open a text file (Python intro)
- ▶ tokenise and normalise text
- ▶ build some text representations



<https://twitter.com/abhi1thakur/status/1575400771541155842>

Representations Revisited

Representations Revisited

1. Use NLTK¹ or Spacy² to tokenize
2. Use `.lower()` to casefold (ignore capitalisation)
3. Use Porter's stemmer to drop suffixes
or use a lemmatiser to find the *actual* root of words
4. Discard stopwords from the text*
5. Build a vectorial representation*

¹<https://www.nltk.org/>

²<https://spacy.io>

Stopwords

Common words in a language that occur with a high frequency, but carry much less substantive information about the meaning of a phrase (Lane et al., 2019, p. 51–54)

Alternative 1 Consider the most frequent tokens in a reference corpus as stopwords (remember Genesis from P4P?)

Alternative 2 Take an existing list of stopwords³

en	es	it
i	a	altri
me	ahora	certa
my	alli	della
it	cerca	nessuna
is	el	prima
do	es	quello
the	unas	solito
will	vez	va
other	yo	via

³For instance, from NLTK, sklearn, or
<https://github.com/stopwords-iso>

Stopwords

Discarding stopwords

- ▶ They are the most frequent tokens in the documents
- ▶ Discarding them reduces the computational effort significantly
- ▶ Typical size of a stopwords list: a few hundred words
- ▶ For some applications (e.g., **topic clustering**), they can be safely discarded
- ▶ For some others (e.g., **dialogue**) they cannot

Stopwords have to be considered with a grain of salt
(as everything in NLP)

More Basic Algebra

Vector representation

BoW

- ▶ A text is represented as the bag (set) of its words
- ▶ It disregards grammar
- ▶ It disregards word order
- ▶ It (can) consider frequency

From (Lane et al., 2019, p. 41)

Dot product

Algebraically, it is the sum of the products of the corresponding entries of the two sequences of numbers $a \cdot b$

$$\begin{aligned} a \cdot b &= \sum_{i=1}^n a_i b_i \\ &= a_1 b_1 + a_2 b_2 + a_3 b_3 + \dots + a_n b_n \end{aligned}$$

```
a = [1, 2, 3]
b = [3, 4, 6]
my_sum = 0
for i in range(len(a)):
    my_sum += a[i] * b[i]
```

There are better —more efficient— ways to compute the dot product!

Now, we can use the dot product to compare two documents (~ similarity)

Vector space model

“[...] an **algebraic** model for representing text documents (and any objects, in general) as vectors of identifiers [...]”⁴

Some applications

- ▶ Relevance rankings in keyword-based search
- ▶ Document clustering to “discover” structure and relations in a text collection

(not the SOTA for most tasks, but it's a *minimum viable product*)

</> **Let us see it working**

⁴https://en.wikipedia.org/wiki/Vector_space_model

Sentiment Analysis

Sentiment Analysis

It **does not** refer to actual sentiment (e.g., love or hate)⁵

It is about **positive** and **negative** (and **neutral**)

a

This monitor is definitely a good value. Does it have superb color and contrast? No. Does it boast the best refresh rate on the market? No. But if you're tight on money, this thing looks and preforms great for the money. It has a Matte screen which does a great job at eliminating glare. The chassis it's enclosed within is absolutely stunning.



POSITIVE

🐦

His [ssa] didnt concede until July 12, 2016. Because he was throwing a tantrum. I can't say this enough: [kcuF] Bernie Sanders.

NEGATIVE

From (Lane et al., 2019, p. 62–65)

⁵That's emotion analysis;
e.g.,  Fernicola et al. (2020);  Zhang et al. (2022)

Valence Aware Dictionary for sEntiment Reasoning
(Hutto and Gilbert, 2014)⁶

- ▶ It has a lexicon packed with tokens and their associated “sentiment” score
- ▶ It counts all tokens belonging to each category: [pos, neu, neg]
... and combine them to determine the sentiment

</> **Let us see it working**

⁶<http://comp.social.gatech.edu/papers/icwsm14.vader.hutto.pdf>
<https://github.com/cjhutto/vaderSentiment>

Tomorrow...

Statistical NLP

References

- Fernicola, F., S. Zhang, F. Garcea, P. Bonora, and A. Barrón-Cedeño
2020. Ariemozione: Identifying emotions in opera verses. In *Italian Conference on Computational Linguistics*.
- Hutto, C. and E. Gilbert
2014. VADER: A parsimonious rule-based model for sentiment analysis of social media text. In *Eighth International Conference on Weblogs and Social Media (ICWSM-14)*, Ann Arbor, MI.
- Lane, H., C. Howard, and H. Hapkem
2019. *Natural Language Processing in Action*. Shelter Island, NY: Manning Publication Co.
- Zhang, S., F. Fernicola, F. Garcea, P. Bonora, and A. Barrón-Cedeño
2022. AriEmozione 2.0: Identifying Emotions in Opera Verses and Arias. *IJCoL*, 8(2).