

92586 Computational Linguistics

Lesson 4. Vector Space Model¹

Alberto Barrón-Cedeño

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

02/03/2022



¹Only a notebook was used in Lesson 3.

Table of Contents

Current Status

Representations Revisited

Sentiment Analysis

Current Status

Current Status

You know...

- ▶ what is computational linguistics / 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

Wondering about your final project?

Get inspiration from last year's
[albarron.github.io/teaching/
computational-linguistics2020/](https://albarron.github.io/teaching/computational-linguistics2020/)

Why not starting with a *fantastic* L^AT_EX template?



github.com/TinfFoil/learning_dit_coli_projecttemplate

Representations Revisited

Representations Revisited

1. Use NLTK² to tokenize
2. Use `.lower()` to ignore capitalisation
3. Use Porter's stemmer to drop suffixes
4. Or use a lemmatiser to find the root of the words
5. Discard stopwords from the text
6. Build a vectorial representation

²<http://www.nltk.org/>

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)

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 $a \cdot b$

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

```
v1 = [1,2,3]
v2 = [3,4,6]
my_sum = 0
for i in range(len(v1)):
    my_sum += v1[i] * v2[i]
```

(there are better —more efficient— ways to compute the dot product)

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
- ▶ Text clustering to “discover” structure and relations in a text collection
- ▶ Reading recommendations

(Not the SoA for most tasks, but it’s a starting point)

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

Sentiment Analysis

Sentiment Analysis

It **does not** refer to real sentiment, such as love or hate⁵
It is about **positive** and **negative** (and **neutral**)



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; see for instance Fernicola et al. 2020

Sentiment Analysis

VADER a rule-based approach

Valence Aware Dictionary for sEntiment Reasoning⁶

- ▶ It has a lexicon packed with tokens and their associated "sentiment" score
- ▶ It counts all tokens belonging to each category: [positive, neutral, negative]

⁶<http://comp.social.gatech.edu/papers/icwsm14.vader.hutto.pdf>

Valence Aware Dictionary for sEntiment Reasoning⁷

- ▶ 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**

⁷<http://comp.social.gatech.edu/papers/icwsm14.vader.hutto.pdf>

Coming soon...

Statistical NLP

References

Lane, H., C. Howard, and H. Hapkem
2019. *Natural Language Processing in Action*. Shelter Island,
NY: Manning Publication Co.