

Some notes on tf-idf

- tf-idf is an effective method for term weighting
- It was originally used for information retrieval, where it brought substantial improvements over other methods
- It is also very effective on text classification when using linear models
- There are some alternatives (e.g., BM25), and many variations: frequencies for TF, or use the log TF
- It has been difficult to improve over it (since 1970's)

A document is more than a BoW

The example document for sentiment analysis

It's a **good** thing most animated sci-fi movies come from Japan, because "titan a.e." is proof that Hollywood **doesn't have a clue** how to do it. I don't know what this film is **supposed** to be about.

- So far, we considered documents as simple BoW words
- BoW representations is surprisingly successful in many fields (IR, spam detection, ...)
- However, word order matters
 - According to a sentiment dictionary, our example contains one **positive** and one **negative** word

Bag of n-grams

- Using n-grams rather than words allows us to capture more information in the data
- We can still use the same weighting methods (tf-idf)
- It is a common practice to use a range of (overlapping) n-grams
- This results in large set of features (millions for most practical applications)
- Data sparsity is a problem for higher order n-grams

	d_1	d_2	d_3	...	d_m
cat	0	3	1	...	4
dog	0	0	3	...	3
book	4	1	4	...	5
...

- Rows of the matrix represent words: words that appear in the same set of documents will be similar to each other
- The columns represent documents: documents with overlapping sets of words will be similar to each other
- Terms do not have to be words, any sequence we can count can be a term
- Contexts do not have to be documents, and any meaningful context can be used instead of documents
- Data is highly correlated (lots of redundancy)
- For practical applications we need huge (but sparse) matrices

A toy example

A four-sentence corpus with *bag of words* (BOW) model.

The corpus:

S1: She likes cats and dogs
S2: He likes dogs and cats
S3: She likes books
S4: He reads books

Term-term (left-context) matrix									
	\emptyset	she	he	likes	reads	cats	dogs	books	and
she	2	0	0	0	0	0	0	0	0
he	2	0	0	0	0	0	0	0	0
likes	0	2	1	0	0	0	0	0	0
reads	0	0	1	0	0	0	0	0	0
cats	0	0	0	1	0	0	0	0	1
dogs	0	0	0	1	0	0	0	0	1
books	0	0	0	1	1	0	0	0	0
and	0	0	0	0	0	1	1	0	0

A simple example from phonetics/phonology

	Vowel	High	Back	Round	Voice	Labial	Nasal	...
i	1	1	0	0	1	0	0	...
a	1	0	1	0	1	0	0	...
m	0	0	0	0	1	0	1	...
n	0	0	0	0	1	1	0	...
p	0	0	0	0	0	0	0	...
b	0	0	0	0	1	0	0	...

- Compared to one-hot representations, these type of representations help identifying similar units

Pointwise mutual information

for term weighting

- Another common weighting method is pointwise mutual information

$$PMI(t, d) = \log \frac{P(t, d)}{P(t)P(d)}$$

- Besides normalizing for 'term frequency/probability', PMI also takes the 'document probability' into account
- Note that 'document' does not have to be a document, any definition of 'context' may result in useful representations (depending on the task)

A document is more than a BoW

The example document for sentiment analysis

It's a **good** thing most animated sci-fi movies come from Japan, because "titan a.e." is proof that Hollywood **doesn't have a clue** how to do it. I don't know what this film is **supposed** to be about.

- So far, we considered documents as simple BoW words
- BoW representations is surprisingly successful in many fields (IR, spam detection, ...)
- However, word order matters
 - According to a sentiment dictionary, our example contains one **positive** and one **negative** word
- Paying attention to longer sequences allows us to get better results

The unreasonable effectiveness of character n-grams

An example document

It's a good thing ...

feature	value
1t	2
t'	1
'g	2
it_1	3
it_2	4
it_3	5
it_4	2
$\text{it}'_1 \text{it}_2$	1
$\text{it}'_2 \text{it}_3$	1
$\text{it}'_3 \text{it}_4$	1
$\text{it}_1 \text{it}_2 \text{it}_3$	1
$\text{it}_1 \text{it}_2 \text{it}_3 \text{it}_4$	2
$\text{it}_2 \text{it}_3 \text{it}_4$	2
...	...

- For a number of text classification tasks (authorship attribution, language detection), character n-gram features found to be effective
- The idea is to use a range of character n-grams

A toy example

A four-sentence corpus with *bag of words* (BOW) model.

Term-document (sentence) matrix

The corpus:

S1: She likes cats and dogs
S2: He likes dogs and cats
S3: She likes books
S4: He reads books

	S1	S2	S3	S4
she	1	0	1	0
he	1	0	1	0
likes	1	1	1	0
reads	0	0	0	1
cats	1	1	0	0
dogs	1	1	0	0
books	0	0	1	1
and	1	1	0	0

What about the linguistic features?

- Linguistically-informed representations is one potential area where linguistics can help building NLP system
- For text classification, the use of 'lexicons' has been common
- Other linguistic features such as
 - lemmas
 - sequences of POS tags
 - parser output: dependency triplets, or partial trees
- are also used in some tasks

Are linguistic features useful at all?

- It is often difficult to get improvements over simple features
- It also makes systems more complex and language dependent
- Linguistic features can particularly be useful if the amount of data is limited
- They are particularly interesting for interpretable and explainable ML/NLP
- Considering the types of linguistic features that help is useful for structuring your input

Final remarks

- Representation of inputs to a ML model is important
- More meaningful/useful representations are likely to improve the systems
- Modern ML methods learn these representations from the data
- Informed/clever ways to represent the data may still be important in some cases (e.g., low-resource scenarios)

Next:

Fri/Mon Learning representations

Some sources of information

- * Jurafsky and Martin (Chapter 6, 2025)

Jurafsky, Daniel and James H. Martin (2025). *Speech and Language Processing: An Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition with Language Models*. 3rd. Online manuscript released January 12, 2025. URL: <https://web.stanford.edu/~jurafsky/slp3/>.