# Big Data & Automated Content Analysis
# Week 6 – Wednesday: »Text as data«

Damian Trilling

d.c.trilling@uva.nl
@damian0604
www.damiantrilling.net

3 March 2021

Afdeling Communicatiewetenschap
Universiteit van Amsterdam

## Today

The bag-of-words (BOW) model

    General idea

Getting a clean BOW representation

    Better tokenization

    Stopword removal

    Pruning

    Stemming and lemmatization

The order of preprocessing steps

How further?

How did it go? And some first feedback on the exam.

The bag-of-words (BOW) model
oooooooooooo

Getting a clean BOW representation
ooooooooooooooooooo

The order of preprocessing steps
ooooo

How further?
ooooo

# Everything clear from last week?

# The bag-of-words (BOW) model

# The bag-of-words (BOW) model

General idea

## A text as a collections of word

Let us represent a string

```
1  t = "This this is is is a test test test"
```

like this:

```
1  from collections import Counter
2  print(Counter(t.split()))
```

```
1  Counter({'is': 3, 'test': 3, 'This': 1, 'this': 1, 'a': 1})
```

Compared to the original string, this representation

- is less repetitive
- preserves word frequencies
- but does *not* preserve word order
- can be interpreted as a vector to calculate with (!!!)

Of course, still a lot of stuff to fine-tune    (for example, This/this)

3

## A text as a collections of word

Let us represent a string

```
1   t = "This this is is is a test test test"
```

like this:
```
1   from collections import Counter
2   print(Counter(t.split()))
```

```
1   Counter({'is': 3, 'test': 3, 'This': 1, 'this': 1, 'a': 1})
```

Compared to the original string, this representation

- is less repetitive
- preserves word frequencies
- but does *not* preserve word order
- can be interpreted as a vector to calculate with (!!!)

*Of course, still a lot of stuff to fine-tune. . . (for example, This/this)*                                         3

## From vector to matrix

If we do this for multiple texts, we can arrange the vectors in a table.

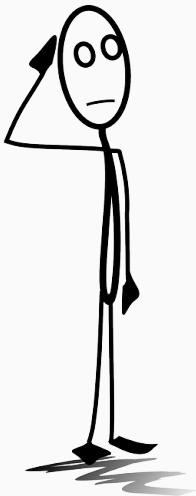t1 = "This this is is is a test test test"
t2 = "This is an example"

|     | a | an | example | is | this | This | test |
|-----|---|----|---------|----|----|------|------|
| t1  | 1 | 0  | 0       | 3  | 1    | 1    | 3    |
| t2  | 0 | 1  | 1       | 1  | 0    | 1    | 0    |

4

*What can you do with such a matrix? Why would you want to represent a collection of texts in such a way?*

**The bag-of-words (BOW) model**
○○○○○○●○○○○○

Getting a clean BOW representation
○○○○○○○○○○○○○○○○○○

The order of preprocessing steps
○○○○○

How further?
○○○○○

# The cell entries: raw counts versus tf·idf scores

- In the example, we entered simple counts (the "term frequency")

*But are all terms equally important?*

## The cell entries: raw counts versus tf·idf scores

- In the example, we entered simple counts (the "term frequency")
- But does a word that occurs in almost all documents contain much information?
- And isn't the presence of a word that occurs in very few documents a pretty strong hint?
- Solution: Weigh by *the number of documents in which the term occurs at least once) (the "document frequency")*

⇒ we multiply the "term frequency" (tf) by the inverse document frequency (idf)

(usually with some additional logarithmic transformation and normalization applied, see
https://scikit-learn.org/stable/modules/generated/sklearn.feature_extraction.text.
TfidfTransformer.html)

8

## The cell entries: raw counts versus tf·idf scores

- In the example, we entered simple counts (the "term frequency")
- But does a word that occurs in almost all documents contain much information?
- And isn't the presence of a word that occurs in very few documents a pretty strong hint?
- **Solution:** Weigh by *the number of documents in which the term occurs at least once) (the "document frequency")*

⇒ we multiply the "term frequency" (tf) by the inverse document frequency (idf)

(usually with some additional logarithmic transformation and normalization applied, see
https://scikit-learn.org/stable/modules/generated/sklearn.feature_extraction.text.
TfidfTransformer.html)

8

## The cell entries: raw counts versus tf·idf scores

- In the example, we entered simple counts (the "term frequency")
- But does a word that occurs in almost all documents contain much information?
- And isn't the presence of a word that occurs in very few documents a pretty strong hint?
- **Solution: Weigh by** *the number of documents in which the term occurs at least once) (the "document frequency")*

$\Rightarrow$ we multiply the "term frequency" (tf) by the inverse document frequency (idf)

(usually with some additional logarithmic transformation and normalization applied, see
https://scikit-learn.org/stable/modules/generated/sklearn.feature_extraction.text.
TfidfTransformer.html)

8

## Is tf·idf always better?

It depends.

- Ultimately, it's an empirical question which works better ($\rightarrow$ weeks on machine learning)

- In many scenarios, "discounting" too frequent words and "boosting" rare words makes a lot of sense (most frequent words in a text can be highly un-informative)

- Beauty of raw tf counts, though: interpretability + describes document in itself, not in relation to other documents

## Internal representations

### Sparse vs dense matrices

- Most are not *not* contained in a given document
- $\rightarrow$ tens of thousands of columns (terms), and one row per document
- Filling all cells is inefficient *and* can make the matrix too large to fit in memory (!!!)
- Solution: store only non-zero values with their coordinates! (sparse matrix)
- dense matrix (or dataframes) not advisable, only for toy examples

## Internal representations

### Little over-generalizing R vs Python remark

Among R users, it is very common to manually inspect document-term matrices, and many operations are done directly on them. In Python, they are more commonly seen as a means to an end (mostly, as input for machine learning).

Many R modules convert to dense matrices: really problematic for larger datasets!

# Getting a clean BOW representation

## Room for improvement

**tokenization**   How do we (best) split a sentence into tokens (terms, words)?

**pruning**   How can we remove unneccessary words?

**lemmatization**   How can we make sure that slight variations of the same word are not counted differently?

# Getting a clean BOW representation

Better tokenization

## OK, good enough, perfect?

### .split()

- space $\rightarrow$ new word
- no further processing whatsoever
- thus, only works well if we do a preprocessing outselves (e.g., remove punctuation)

```
1  docs = ["This is a text", "I haven't seen John's derring-do. Second
        sentence!"]
2  tokens = [d.split() for d in docs]
```

```
1  [['This', 'is', 'a', 'text'], ['I', "haven't", 'seen', "John's", 'derring-do.', 'Second', '
        sentence!']]
```

## OK, good enough, perfect?

### Tokenizers from the NLTK pacakge

- multiple improved tokenizers that can be used instead of .split()
- e.g., Treebank tokenizer:
  - split standard contractions ("don't")
  - deals with punctuation

```
1  from nltk.tokenize import TreebankWordTokenizer
2  tokens = [TreebankWordTokenizer().tokenize(d) for d in docs]
```

```
1  [['This', 'is', 'a', 'text'], ['I', 'have', 'n't', 'seen', 'John', "'s", 'derring-do.', 'Second
         ', 'sentence', '!']]
```

Notice the failure to split the . at the end of the first sentence in the second doc. That's because TreebankWordTokenizer expects *sentences* as input. See book for a solution.

The bag-of-words (BOW) model    **Getting a clean BOW representation**    The order of preprocessing steps    How further?

00000000000    0000000●00000000000    00000    00000

OK, so we can tokenize with a list comprehension (and that's often a good idea!). But what if we want to *directly* get a DTM instead of lists of tokens?

## OK, good enough, perfect?

### scikit-learn's CountVectorizer (default settings)

- applies lowercasing

- deals with punctuation etc. itself

- minimum word length $> 1$

- more technically, tokenizes using this regular expression:
  `r"(?u)\b\w\w+\b"` [1]

```
1  from sklearn.feature_extraction.text import CountVectorizer
2  cv = CountVectorizer()
3  dtm_sparse = cv.fit_transform(docs)
```

_____

[1] ?u = support unicode, \b = word boundary

15

## OK, good enough, perfect?

### CountVectorizer supports more

- stopword removal

- custom regular expression

- or even using an external tokenizer

- ngrams instead of unigrams

see https://scikit-learn.org/stable/modules/generated/sklearn.feature_extraction.text.
CountVectorizer.html

### Best of both worlds

Use the Count vectorizer with a NLTK-based external tokenizer! (see book)

16

The bag-of-words (BOW) model    **Getting a clean BOW representation**    The order of preprocessing steps    How further?

00000000000     00000000000000     00000     00000

# OK, good enough, perfect?

## CountVectorizer supports more

- stopword removal

- custom regular expression

- or even using an external tokenizer

- ngrams instead of unigrams

see https://scikit-learn.org/stable/modules/generated/sklearn.feature_extraction.text.
CountVectorizer.html

## Best of both worlds

Use the Count vectorizer with a NLTK-based external tokenizer! (see book)

# Getting a clean BOW representation

Stopword removal

## Stopword removal

### What are stopwords?

- Very frequent words with little inherent meaning
- the, a, he, she, ...
- context-dependent: if you are interested in gender, he and she are no stopwords.
- Many existing lists as basis

When using the CountVectorizer, we can simply provide a stopword list.

But we can also remove stopwords "by hand" (next slide):

## Stopword removal

```
1  from nltk.corpus import stopwords
2  mystopwords = stopwords.words("english")
3  mystopwords.extend(["test", "this"])
4
5  def tokenize_clean(s, stoplist):
6      cleantokens = []
7      for w in TreebankWordTokenizer().tokenize(s):
8          if w.lower() not in stoplist:
9              cleantokens.append(w)
10     return cleantokens
11
12 tokens = [tokenize_clean(d, mystopwords) for d in docs]
```

```
1  [['text'], ["n't", 'seen', 'John', 'derring-do.', 'Second', 'sentence', '!']]
```

### You can do more!

For instance, in line 8, you could add an or statement to also exclude punctuation.

18

# Getting a clean BOW representation

## Pruning

# General idea

- Idea behind both stopword removal and tf·idf: too frequent words are uninformative
- (possible) downside stopword removal: a priori list, does not take empirical frequencies in dataset into account
- (possible) downside tf·idf: does not reduce number of features

Pruning: remove all features (tokens) that occur in less than X or more than X of the documents

## General idea

- Idea behind both stopword removal and tf·idf: too frequent words are uninformative
- (possible) downside stopword removal: a priori list, does not take empirical frequencies in dataset into account
- (possible) downside tf·idf: does not reduce number of features

Pruning: remove all features (tokens) that occur in less than X or more than X of the documents

## General idea

- Idea behind both stopword removal and tf·idf: too frequent words are uninformative
- (possible) downside stopword removal: a priori list, does not take empirical frequencies in dataset into account
- (possible) downside tf·idf: does not reduce number of features

Pruning: remove all features (tokens) that occur in less than X or more than X of the documents

## General idea

- Idea behind both stopword removal and tf·idf: too frequent words are uninformative
- (possible) downside stopword removal: a priori list, does not take empirical frequencies in dataset into account
- (possible) downside tf·idf: does not reduce number of features

Pruning: remove all features (tokens) that occur in less than X or more than X of the documents

## CountVectorizer, only stopword removal

```
1  from sklearn.feature_extraction.text import CountVectorizer,
       TfidfVectorizer
2  myvectorizer = CountVectorizer(stop_words=mystopwords)
```

CountVectorizer, better tokenization, stopword removal (pay attention that stopword list uses same tokenization!):

```
1  myvectorizer = CountVectorizer(tokenizer = TreebankWordTokenizer().
       tokenize, stop_words=mystopwords)
```

Additionally remove words that occur in more than 75% or less than $n = 2$ documents:

```
1  myvectorizer = CountVectorizer(tokenizer = TreebankWordTokenizer().
       tokenize, stop_words=mystopwords, max_df=.75, min_df=2)
```

All togehter: tf·idf, explicit stopword removal, pruning

```
1  myvectorizer = TfidfVectorizer(tokenizer = TreebankWordTokenizer().
       tokenize, stop_words=mystopwords, max_df=.75, min_df=2)
```

*What is "best"? Which (combination of) techniques to use, and how to decide?*

# Getting a clean BOW representation

Stemming and lemmatization

## Stemming and lemmatization

- Stemming: reduce words to its stem by removing last part (drinking → drink)
- Lemmatization: find word that you would need to look up in a dictionary (drinking → drink, but also went → go)
- stemming is simpler than lemmatization
- lemmatization often better

Example below: tokenization and lemmatization with spacy in one go:

```
1  import spacy
2  nlp = spacy.load('en') # potentially you need to install the language
        model first
3  lemmatized_tokens = [[token.lemma_ for token in nlp(doc)] for doc in
        docs]
```

```
1  [['this', 'be', 'a', 'text'], ['-PRON-', 'have', 'not', 'see', 'John', "'s", 'darling', '.', 'do
        ', '.', 'second', 'sentence', '.']]
```

## Stemming and lemmatization

- Stemming: reduce words to its stem by removing last part (drinking → drink)
- Lemmatization: find word that you would need to look up in a dictionary (drinking → drink, but also went → go)
- stemming is simpler than lemmatization
- lemmatization often better

Example below: tokenization and lemmatization with spacy in one go:

```
1  import spacy
2  nlp = spacy.load('en') # potentially you need to install the language
       model first
3  lemmatized_tokens = [[token.lemma_ for token in nlp(doc)] for doc in
       docs]
```

```
1  [['this', 'be', 'a', 'text'], ['-PRON-', 'have', 'not', 'see', 'John', "'s", 'derring', '-', 'do
       ', '.', 'second', 'sentence', '!']]
```

# The order of preprocessing steps

## Option 1

### Preprocessing only through Vectorizer

"Just use CountVectorizer or Tfidfvectorizer with the appropriate options."

- pro: No double work, efficient if your main goal is a sparse matrix (for ML?) anyway

- con: you cannot "see" the preprocessed texts

## Option 2

### Extensive preprocessing without Vectorizer

"Remove stopwords, punctuation etc. and store in a string with spaces"

```
1  cleaneddocs = [" ".join(re.findall(r"\w\w+", d)).lower() for d in docs]
2  cleaneddocswithoutstopwords = [" ".join([w for w in d.split() if w not
       in mystopwords]) for d in cleaneddocs]
```

```
1  ['this is text', 'haven seen john derring do second sentence']
2  ['text', 'seen john derring second sentence']
```

Yes, this list comprehension looks scary – you can make a more elaborate for loop instead

- pro: you can read (and store!) the preprocessed docs
- pro: even the most stupid vectorizer (or wordcloud tool) can split the resulting string later on
- con: potentially double work (for you *and* the computer)

24

How would you do it?

I tend to go for Option 2 because

- I like to inspect a sample of the documents
- I can re-use the cleaned docs irrespective of the Vectorizer

But sometimes, I opt of Option 1 instead because

- I want to systematically compare the effect of different choices in a machine learning pipeline (then I can simply vary the vectorizer instead of the data)
- I want to use techniques that are geared towards little or no preprocessing (deep learning)

# How further?

## Main takeaway

- It matters how you transform your text into numbers ("vectorization").

- Preprocessing matters, be able to make informed choices.

- Keep this in mind when we will discuss Machine Learning! It will come back throughout Part II!

- Once you vectorized your texts, you can do all kinds of calculations (random example: get the cosine similarity between two texts)

27

## More NLP

*n*-**grams** Consider using *n*-grams instead of unigrams

**collocations** *n*grams that appear more frequently than expected

**POS-tagging** grammatical function ("part-of-speach") of tokens

**NER** named entity recognition (persons, organizations, locations)

## More NLP

I **really** recommend looking into spacy (https://spacy.io) for advanced natural language processing, such as part-of-speech-tagging and named entity recogntion.

The bag-of-words (BOW) model
○○○○○○○○○○○

Getting a clean BOW representation
○○○○○○○○○○○○○○○○○

The order of preprocessing steps
○○○○○

How further?
○○○○●

## Friday

### Based on today's lecture and Chapter 10...

Try to take some of the data from last week and

- preprocess them (in different ways)
- vectorize them
- give a (tabular and/or graphical) overview of tokens (unigrams, bigrams, collocations).

Then,

- compare that bottom-up approach with a top-down (keyword or regular-expression based) approach