

Big Data and Automated Content Analysis (6EC)

Week 6: »Processing textual data«

Monday

Anne Kroon

a.c.kroon@uva.nl, @annekroon

May 9, 2022

Today

Basic string operations

Regular expressions

What is a regexp?

Using a regexp in Python

The bag-of-words (BOW) model

General idea

A cleaner BOW representation

Better tokenization

Stopword removal

Pruning

Stemming and lemmatization



Everything clear from last week?

Basic string operations

oooooooo

Regular expressions

o
oooooo
oooooooo

The BOW

o
oooooooooooo
oooooooooooooooo
ooooo
ooooo

Next steps

ooo

This week, we will get a general overview of working with textual data. Combining the knowledge from this week with last week gives you all blocks you need to do cool automated content analyses.

Basic string operations

Working with strings

1. string methods that every string has (`"hello".upper()`)
2. functions that take a string as input (`len("hello")`)
3. pandas column string methods
(`df["somecolumn"].str.upper()`)
4. applying string methods or functions to a pandas column
(`df["somecolumn"].apply(len)` or
`df["somecolumn"].apply(lambda x: x.upper())`)

For today, we assume that our data are a list of strings – adapt accordingly for pandas.

An example says more than 1000 words...

```

1  # probably read from text file(s) instead, you learned that already...
2  data = [ "I <b>really</b> liked this movie! It was great. ", " What
   ↪  an awful movie", "Awesome!!!"]
3
4  data_stripped = [e.strip() for e in data]
5  data_lower = [e.lower() for e in data_stripped]
6  data_clean = [e.replace("<b>", "").replace("</b>", "") for e in
   ↪  data_lower]
7
8  # or, more efficient, in one single step:
9  data_clean2 = [e.strip().lower().replace("<b>", "").replace("</b>", "")
   ↪  for e in data]

```

Two examples says even more:

```

1  from string import punctuation
2
3  # punctuation is just the string '!"#$%&'()*+,-./:;<=>?@[\\]^_`{|}~'
4
5  text = "This is a test! Let's get rid (of) punct&"
6
7  # we make a list of each character in the text but only if it is not
8  # a punctuation sign. The, we join the elements of the list directly
9  # to each other without anything between it ("")
10 cleantext = "".join([c for c in text if c not in punctuation])

```


Combine both

```

1  from string import punctuation
2
3  def strip_punctuation(text):
4      return "".join([c for c in text if c not in punctuation])
5
6  data_clean3 = [strip_punctuation(e).strip().lower()\
7      .replace("<b>","").replace("</b>","") for e in data]

```

The toolbox at a glance

Slicing

`mystring[2:5]` to get the characters with indices 2,3,4

String methods

- `.lower()` returns lowercased string
- `.strip()` returns string without whitespace at beginning and end
- `.find("bla")` returns index of position of substring "bla" or -1 if not found
- `.replace("a","b")` returns string with "a" replaced by "b"
- `.count("bla")` counts how often substring "bla" occurs
- `.isdigit()` true if only numbers

Use tab completion for more!

From test to large-scale analysis: General approach

1. Take a single string and test your idea

```
1 t = "This is a test test test."  
2 print(t.count("test"))
```

2a. You'd assume it to return 3. If so, scale it up:

```
1 results = []  
2 for t in listwithallmytexts:  
3     r = t.count("test")  
4     print(f"{t} contains the substring {r} times")  
5     results.append(r)
```

2b. If you *only* need to get the list of results, a list comprehension is more elegant:

```
1 results = [t.count("test") for t in listwithallmytexts]
```

General approach

Test on a single string, then make a for loop or list comprehension!

Own functions

If it gets more complex, you can write your own function and then use it in the list comprehension:

```

1 def mycleanup(t):
2     # do sth with string t here, create new string t2
3     return t2
4
5 results = [mycleanup(t) for t in allmytexts]
```

Pandas string methods as alternative

If you select column with strings from a pandas dataframe, pandas offers a collection of string methods (via `.str.`) that largely mirror standard Python string methods:

```
1 df['newcoloumnwithresults'] = df['columnwithtext'].str.count("bla")
```

To pandas or not to pandas for text?

Partly a matter of taste.

Not-too-large dataset with a lot of extra columns? Advanced statistical analysis planned? Sounds like pandas.

It's mainly a lot of text? Wanna do some machine learning later on anyway? It's large and (potentially) messy? Doesn't sound like pandas is a good idea.

Regular expressions

Regular expressions

What is a regexp?

Regular Expressions: What and why?

What is a regexp?

- a *very* widespread way to describe patterns in strings
- Think of wildcards like `*` or operators like OR, AND or NOT in search strings: a regexp does the same, but is *much* more powerful
- You can use them in many editors (!), in the Terminal, in STATA ...and in Python

A more powerful tool

An example

- We want to remove everything but words from a tweet
- We can do so by calling the `.replace()` method multiple times (for each unwanted character)
- But we can better do this with a regular expression instead:
[`^a-zA-Z`] matches anything that is not a letter

Basic regexp elements

Alternatives

`[TtFf]` matches either T or t or F or f

`Twitter|Facebook` matches either Twitter or Facebook

`.` matches any character

Repetition

`?` the expression before occurs 0 or 1 times

`*` the expression before occurs 0 or more times

`+` the expression before occurs 1 or more times

regex quizz

Which words would be matched?

1. [Pp]ython
2. [A-Z]+
3. RT ? : ? @ [a-zA-Z0-9]+

Basic string operations
oooooooo

Regular expressions
o
o
oooo●
ooooooo

The BOW
o
oooooooooooo
oooooooooooooooooooo
ooooo
ooooo

Next steps
ooo

What else is possible?

See the table in the book!

Regular expressions

Using a regexp in Python

How to use regular expressions in Python

The module `re`*

`re.findall("[Tt]witter|[Ff]acebook", testo)` returns a list with all occurrences of Twitter or Facebook in the string called `testo`

`re.findall("[0-9]+[a-zA-Z]+", testo)` returns a list with all words that start with one or more numbers followed by one or more letters in the string called `testo`

`re.sub("[Tt]witter|[Ff]acebook", "a social medium", testo)` returns a string in which all occurrences of Twitter or Facebook are replaced by "a social medium"

Use the less-known but more powerful module `regex` instead to support all dialects used in the book

How to use regular expressions in Python

The module re

`re.match(" +([0-9]+) of ([0-9]+) points",line)` returns `None` unless it *exactly* matches the string `line`. If it does, you can access the part between `()` with the `.group()` method.

Example:

```
1 line="                2 of 25 points"
2 result=re.match(" +([0-9]+) of ([0-9]+) points",line)
3 if result:
4     print ("Your points:",result.group(1))
5     print ("Maximum points:",result.group(2))
```

Your points: 2

Maximum points: 25

Possible applications

Data preprocessing

- Remove unwanted characters, words, ...
- Identify *meaningful* bits of text: usernames, headlines, where an article starts, ...
- filter (distinguish relevant from irrelevant cases)

Possible applications

Data analysis: Automated coding

- Actors
- Brands
- links or other markers that follow a regular pattern
- Numbers (!)

Example 1: Counting actors

```
1 import re, csv
2 from glob import glob
3 counts1=[]
4 counts2=[]
5 filenames = glob("/home/damian/articles/*.txt")
6
7 for fn in filenames:
8     with open(fn) as fi:
9         artikel = fi.read()
10        artikel = artikel.replace('\n', ' ')
11
12        ↪ counts1.append(len(re.findall('Israel.*(minister|politician.*|[Aa]utl
13        counts2.append(len(re.findall('[Pp]alest',artikel)))
14
15 output=zip(filenames, counts1, counts2)
16 with open("results.csv", mode='w',encoding="utf-8") as fo:
17     writer = csv.writer(fo)
18     writer.writerows(output)
```

Example 2: Parsing semi-structured data

If your data look like this, you can loop over the lines and use regular expressions to extract the info you need!

```
1           All Rights Reserved
2
3           2 of 200 DOCUMENTS
4
5           De Telegraaf
6
7           21 maart 2014 vrijdag
8
9 Brussel bereikt akkoord aanpak probleembanken;
10 ECB krijgt meer in melk te brokkelen
11
12 SECTION: Finance; Blz. 24
13 LENGTH: 660 woorden
14
15 BRUSSEL Europa heeft gisteren op de valreep een akkoord bereikt
16 over een saneringsfonds voor banken. Daarmee staat de laatste
```

Practice yourself!

Take some time to write some regular expressions. Write a script that

- extracts URLs from a list of strings
- removes everything that is not a letter or number from a list of strings

(first develop it for a single string, then scale up)

More tips: <http://www.pyregex.com/>

The BOW

The BOW

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

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



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

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)

Is tf·idf always better?

It depends.

- Ultimately, it's an empirical question which works better (→ 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
- → 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!

The BOW

A cleaner BOW representation

Room for improvement

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

pruning How can we remove unnecessary words?

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

OK, good enough, perfect?

.split()

- space → new word
- no further processing whatsoever
- thus, only works well if we do a preprocessing ourselves (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 package

- multiple improved tokenizers that can be used instead of `.split()`
- e.g., Treebank tokenizer:
 - split standard contractions ("don't")
 - deals with punctuation
 - BUT: Assumes lists of *sentences*.
- Solution: Build an own (combined) tokenizer (next slide)!

OK, good enough, perfect?

```
1 import nltk
2 import regex
3
4 class MyTokenizer:
5     def tokenize(self, text):
6         tokenizer = nltk.tokenize.TreebankWordTokenizer()
7         result = []
8         word = r"\p{letter}"
9         for sent in nltk.sent_tokenize(text):
10             tokens = tokenizer.tokenize(sent)
11             tokens = [t for t in tokens
12                      if regex.search(word, t)]
13             result += tokens
14         return result
15
16 mytokenizer = MyTokenizer()
17 tokens = [mytokenizer.tokenize(d) for d in docs]
```

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



Can you (try to) explain the code?

Basic string operations
○○○○○○○○○

Regular expressions
○
○○○○○○○
○○○○○○○○○

The BOW
○
○○○○○○○○○○○
○○○○○○○●○○○○○○○
○○○○○
○○○○○

Next steps
○○○

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

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 the custom NLTK-based external tokenizer we created before! `cv = CountVectorizer(tokenizer=mytokenizer.tokenize)`

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 stopwords list.

But we can also remove stopwords “by hand” of course using either a for loop (like we did for punctuation removal) or by modifying the tokenizer (try it!).

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 stopwords removal

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

CountVectorizer, other tokenization, stopwords removal (pay attention that stopwords 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 together: tf-idf, explicit stopwords 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?

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]
```

The BOW

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



How would you do it?

Sometimes, I 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 at other times, 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)

The BOW

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)

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-speech”) 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 recognition.



Any questions?

Next steps

Basic string operations
○○○○○○○○○

Regular expressions
○
○○○○○
○○○○○○○

The BOW
○
○○○○○○○○○
○○○○○○○○○○○○○
○○○○○
○○○○○

Next steps
○●○

Make sure you understood all of today's
concepts.

Re-read the chapters.

I prepared exercises to work on *during* the
Thursday meeting (alone or in teams):

[https://github.com/uvacw/teaching-bdaca/blob/
main/6ec-course/week06/exercises/](https://github.com/uvacw/teaching-bdaca/blob/main/6ec-course/week06/exercises/)

