

# Big Data & Automated Content Analysis

## Week 5 – Wednesday: »Working with text«

---

Anne Kroon

a.c.kroon@uva.nl

@annekroon

29 April 2021

Afdeling Communicatiewetenschap  
Universiteit van Amsterdam

# Today

Bottom-up vs. top-down

Approaches to working with text

The toolbox

From test to large-scale analysis

ACA using regular expressions

What is a regexp?

Using a regexp in Python

Bottom-up vs. top-down

ooooooo

Approaches to working with text

o

oo

ooo

oooo

Regular expressions

oo

oooooo

oooooooooooo

Everything clear from last week?

## Bottom-up vs. top-down

---

Bottom-up vs. top-down

●○○○○○

Approaches to working with text

○  
○  
○○  
○○○

Regular expressions

○○  
○○○○○  
○○○○○○○○○○

Automated content analysis can be either **bottom-up** (inductive, explorative, pattern recognition, . . . ) or **top-down** (deductive, based on a-priori developed rules, . . . ). Or in between.

# The ACA toolbox

	Methodological approach		
	<i>Counting and Dictionary</i>	<i>Supervised Machine Learning</i>	<i>Unsupervised Machine Learning</i>
<b>Typical research interests and content features</b>	visibility analysis sentiment analysis subjectivity analysis	frames topics gender bias	frames topics
<b>Common statistical procedures</b>	string comparisons counting	support vector machines naive Bayes	principal component analysis cluster analysis latent dirichlet allocation semantic network analysis

deductive

inductive

Boumans2016

## Bottom-up vs. top-down

### Bottom-up

- Count most frequently occurring words
- Maybe better: Count combinations of words  $\Rightarrow$  Which words co-occur together?

We *don't* specify what to look for in advance

### Top-down

- Count frequencies of pre-defined words
- Maybe better: patterns instead of words

We *do* specify what to look for in advance

# A simple bottom-up approach

```
1 from collections import Counter
2
3 texts = ["I really really really love him, I do", "I hate him"]
4
5 for t in texts:
6     print(Counter(t.split()).most_common(3))
```

```
1 [('really', 3), ('I', 2), ('love', 1)]
2 [('I', 1), ('hate', 1), ('him', 1)]
```



## A simple top-down approach

```
1 texts = ["I really really really love him, I do", "I hate him"]
2 features = ['really', 'love', 'hate']
3
4 for t in texts:
5     print(f"\nAnalyzing '{t}':")
6     for f in features:
7         print(f"{f} occurs {t.count(f)} times")
```

```
1 Analyzing 'I really really really love him, I do':
2 really occurs 3 times
3 love occurs 1 times
4 hate occurs 0 times
5
6 Analyzing 'I hate him':
7 really occurs 0 times
8 love occurs 0 times
9 hate occurs 1 times
```



*When would you use which approach?*

## Approaches to working with text

---

# Approaches to working with text

---

## The toolbox

# The toolbox

## 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 where "a" is replaced by "b"
- `.count("bla")` counts how often substring "bla" occurs

# The toolbox

Regular expressions

(today)

# Approaches to working with text

---

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

---

Automated content analysis using 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

## An example

### Let's say...

- We wanted to remove everything but words from a tweet
- We could do so by calling the `.replace()` method
- We could do this with a regular expression as well:  
    `[^a-zA-Z]` would match 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 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] \*

## 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 count1_list=[]
4 count2_list=[]
5 filename_list = glob("/home/damian/articles/*.txt")
6
7 for fn in filename_list:
8     with open(fn) as fi:
9         artikel = fi.read()
10        artikel = artikel.replace('\n',' ')
11
12        count1 = len(re.findall('Israel.*(minister|politician.*|[Aa]
13                               uthorit)',artikel))
14
15        count2 = len(re.findall('[Pp]alest',artikel))
16
17        count1_list.append(count1)
18        count2_list.append(count2)
19
20 output=zip(filename_list,count1_list, count2_list)
21 with open("results.csv", mode='w',encoding="utf-8") as fo:
22     writer = csv.writer(fo)
23     writer.writerows(output)
```



## Example 2: Which number has this Lexis Nexis article?

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 probleebanken;  
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

## Example 2: Check the number of a lexis nexis article

```
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 probleebanken;
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
```

```
1 for line in tekst:
2     matchObj=re.match(r" +([0-9]+) of ([0-9]+) DOCUMENTS",line)
3     if matchObj:
```

## Practice yourself!

Let's 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/>

## Next meetings

### Practice yourself...

Write your own ACA script! see <https://github.com/annekroon/bdaca-6ec/tree/master/6ec/week05/exercises/exercise.md>.

### TAKE HOME EXAM

Any questions? Deadline: Monday, 23.59

