

# Big Data and Automated Content Analysis

## Week 4 – Thursday

### » Sentiment Analysis «

Anne Kroon

a.c.kroon@uva.nl  
@annekroon

Afdeling Communicatiewetenschap  
Universiteit van Amsterdam

April 25, 2019

# Today

- ① API assignment last week
- ② Different types of analysis
  - What can we do?
  - Systematizing analytical approaches
- ③ Data analysis 1: Sentiment analysis
  - What is it?
  - Bag-of-words approaches
  - Advanced approaches
  - A sentiment analysis tailored to your needs!
  - Packages for sentiment analysis
  - A recipe
  - Machine Learning as alternative
- ④ Take-home message, next meetings, & exam

## Creating an URL to make an API request

```
1 import requests
2 import json
3
4 base_url = 'https://www.rijksmuseum.nl/api/pages/en/{}?key={}&format=
    json'
5 page_of_interest = 'whats-on/exhibitions-past'
6 key = 'YOURKEY'
7
8 full_url = base_url.format(page_of_interest, key)
```

NB: the built-in `.format()` method returns a formatted value / representation. More specifically, it replaces the `'{'` in a string by the argument specified between `'()'`

NB: assign your unique API-key requested to the variable `key`

## Making an API request

```
1 r = requests.get(full_url)
2 print('Response HTTP Status Code: {}'.format(r.status_code))
3
4 data = json.loads(r.content.decode('utf-8'))
```

NB: Response HTTP Status Code: 200 means success



## Now its your turn!

## How to get started?

Search and consult online documentation on using APIs (e.g., on the organization's webpage, github page).

Consult Stackoverflow or other online sources for help.



# Reddit API

```
1 reddit_dict = { "title": [], "score": [], "id": [], "url": [], "comms_num":  
    [], "created": [], "body": []}  
2  
3 for submission in reddit.subreddit('learnprogramming').hot(limit=15):  
4     reddit_dict["title"].append(submission.title)  
5     reddit_dict["score"].append(submission.score)  
6     reddit_dict["id"].append(submission.id)  
7     reddit_dict["url"].append(submission.url)  
8     reddit_dict["comms_num"].append(submission.num_comments)  
9     reddit_dict["created"].append(submission.created)  
10    reddit_dict["body"].append(submission.selftext)
```

`https://praw.readthedocs.io/en/latest/code_overview/models/  
subreddit.html`





**What do you think? What are interesting methods to analyze large data sets (like, e.g., social media data? What questions can they answer?)**

# What else can we do?

## For example

- sentiment analysis
- automated coding with regular expressions
- natural language processing
- supervised and unsupervised machine learning
- network analysis



## Systematizing analytical approaches

Taking the example of Twitter:

## Analyzing the *structure*

- Number of Tweets over time
- singleton/retweet ratio
- Distribution of number of Tweets per user
- Interaction networks

⇒ **Focus on the amount of content and on the question who interacts with whom, not on what is said**

Bruns, A., & Stieglitz, S. (2013). Toward more systematic Twitter analysis: metrics for tweeting activities. *International Journal of Social Research Methodology*. doi:10.1080/13645579.2012.756095







# Automated Content Analysis

	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 naïve Bayes	principal component analysis cluster analysis latent dirichlet allocation semantic network analysis

deductive
inductive

Boumans, J.W., & Trilling, D. (2016). Taking stock of the toolkit: An overview of relevant automated content analysis approaches and techniques for digital journalism scholars. *Digital Journalism*, 4, 1. 8–23.

## Sentiment analysis

- \* Less sophisticated approaches do not see this as a separate dimension but simply calculate  $objectivity = 1 - (negativity + positivity)$



## Who uses it?

- Companies
- especially for Web Analytics
- Social Scientists
- applications in data journalism, politics, ...

Many references to examples in Mostafa (2013).

⇒ Cases in which you have a huge amount of data or real-time data and you want to get an idea of the tone.

Mostafa, M. M. (2013). More than words: Social networks' text mining for consumer brand sentiments. *Expert Systems with Applications*, 40(10), 4241– 4251. doi:10.1016/j.eswa.2013.01.019

## Data analysis 1: Sentiment analysis

### Bag-of-words approaches

## Bag-of-words approaches

## How does it work?

- We take each word of a text and look if it's positive or negative.
  - Most simple way: compare it with a list of negative words and with a list of positive words (That's what Mostafa (2013) did)
  - More advanced: look up a subjectivity score from a table
- e.g., add up the scores and average them.

If you were to run an analysis like the one by Mostafa (2013), how could you do this?



# How to do this

(given a *string* `tekst` that you want to analyze and two *lists* of strings with negative and positive words, `lijstpos=["great","fantastic",...,"perfect"]` and `lijstneg`)

```

1 sentiment=0
2 for woord in tekst.split():
3     if woord in lijstpos:
4         sentiment=sentiment+1 #same as sentiment+=1
5     elif woord in lijstneg:
6         sentiment=sentiment-1 #same as sentiment-=1
7 print (sentiment)

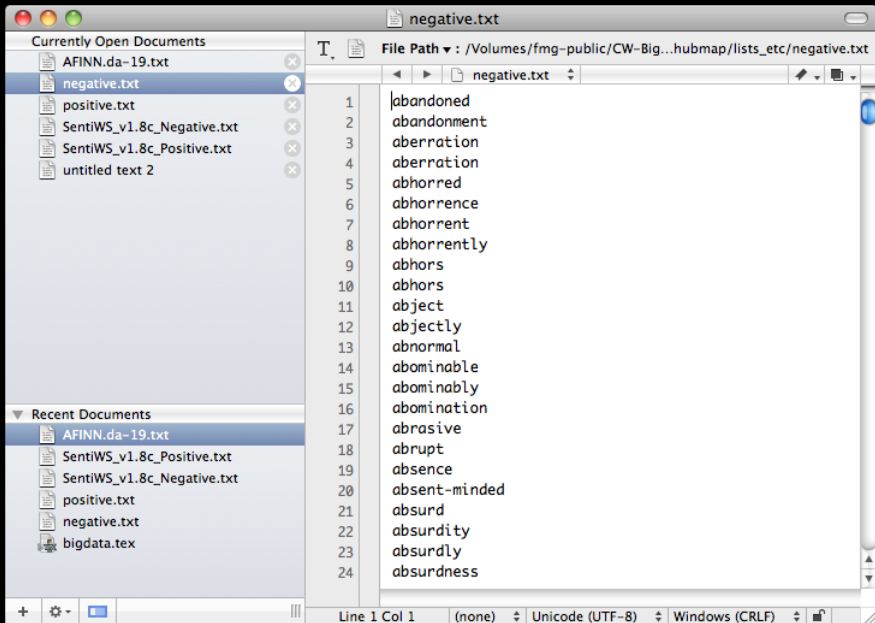
```

# Do we need to have the lists in our program itself?

No.

You could have them in a separate text file, one per row, and then read that file directly to a list.

```
1 poslijst=open("filewithonepositivewordperline.txt").read().splitlines()
2 neglijst=open("filewithonenegativewordperline.txt").read().splitlines()
```



## More advanced versions

- CSV files or similar tables with weights
- Or some kind of dict?

AFINN.da-19.txt

Currently Open Documents

- AFINN.da-19.txt
- negative.txt
- positive.txt
- SentiWS\_v1.8c\_Negative.txt
- SentiWS\_v1.8c\_Positive.txt
- untitled text 2

Recent Documents

- AFINN.da-19.txt
- SentiWS\_v1.8c\_Positive.txt
- SentiWS\_v1.8c\_Negative.txt
- positive.txt
- negative.txt
- bigdata.tex

File Path: /Volumes/fmg-public/CW-Big...ap/lists\_etc/AFINN.da-19.txt

AFINN.da-19.txt

1	absorberet	1
2	acceptere	1
3	accepterede	1
4	accepterer	1
5	accepteres	1
6	accepteret	1
7	advare	-2
8	advarede	-2
9	advarer	-2
10	advaret	-2
11	advarsel	-3
12	advarsler	-3
13	advarslerne	-3
14	afbrudt	-2
15	afbryde	-2
16	afbrydelse	-2
17	afbrydelser	-2
18	afbrydelserne	-2
19	afbryder	-2
20	affald	-1
21	afgift	-1
22	afgifter	-1
23	afhængig	-1
24	afhængige	-1

Line 1 Col 1 (none) Unicode (UTF-8, with BOM) Wind...CRLF

Currently Open Documents

- AFINN.da-19.txt
- negative.txt
- positive.txt
- SentiWS\_v1.8c\_Negative.txt
- SentiWS\_v1.8c\_Positive.txt
- untitled text 2

Recent Documents

- AFINN.da-19.txt
- SentiWS\_v1.8c\_Positive.txt
- SentiWS\_v1.8c\_Negative.txt
- positive.txt
- negative.txt
- bigdata.tex

File Path: /Volumes/fmg-public/CW-Big...tc/SentiWS\_v1.8c\_Negative.txt

SentiWS\_v1.8c\_Negative.txt

1	Abbau INN -0.058	Abbaus,Abbaues,Abbauen,Abbaue
2	Abbruch INN -0.0048	
...	Abbruches,Abbrüche,Abbruchs,Abbrüchen	
3	Abdankung INN -0.0048	Abdankungen
4	Abdämpfung INN -0.0048	Abdämpfungen
5	Abfall INN -0.0048	
...	Abfalles,Abfälle,Abfalls,Abfällen	
6	Abfuhr INN -0.3367	Abfahren
7	Abgrund INN -0.3465	
8	Abhängigkeit INN -0.3653	Abhängigkeiten
9	Ablehnung INN -0.5118	Ablehnungen
10	Ablenkung INN -0.0435	Ablenkungen
11	Abnahme INN -0.0048	Abnahmen
12	Abneigung INN -0.0048	Abneigungen
13	Abnutzung INN -0.0048	
14	Abriss INN -0.0048	
...	Abrisse,Abrissen,Abrisses,Abriss	
15	Abrutsch INN -0.0048	
...	Abrutschen,Abrutsche,Abrutsches,Abrutschs	
16	Abschaffung INN -0.058	Abschaffungen
17	Abschreckung INN -0.0048	Abschreckungen
18	Abschreibung INN -0.3345	Abschreibungen
19	Abschuß INN -0.0048	
20	Abschwächung INN -0.1935	Abschwächungen

Line 1 Col 1 (none) Unicode (UTF-8) Unix (LF) 216...

# Mustafa 2013: Interpreting the output

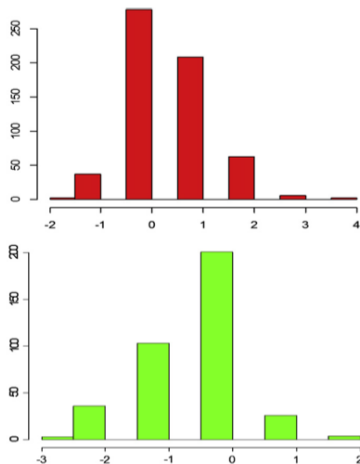


Fig. 5. Sentiment scores for Nokia (top) and Pfizer (bottom). X-axis represents score distributions, Y-axis represents count/frequencies.

Your thoughts?

- each word counts equally (1)
- many tweets contain no words from the list. What does this mean?
- Ways to improve BOW approaches?

# Bag-of-words approaches

## pro

- easy to implement
- easy to modify:
  - add or remove words
  - make new lists for other languages, other categories (than positive/negative), ...
- easy to understand (transparency, reproducibility)

e.g., Schut, L. (2013). Verenigde Staten vs. Verenigd Koninkrijk: Een automatische inhoudsanalyse naar verklarende factoren voor het gebruik van positive campaigning en negative campaigning door vooraanstaande politici en politieke partijen op Twitter. *Bachelor Thesis*, Universiteit van Amsterdam.



# Bag-of-words approaches

## con

- simplistic assumptions
- e.g., intensifiers cannot be interpreted ("really" in "really good" or "really bad")
- or, even more important, negations.

# Data analysis 1: Sentiment analysis

## Advanced approaches

# Improving the BOW approach

## Example: The Sentistrength algorithm

- $-5 \dots -1$  and  $+1 \dots +5$
- spelling correction
- "booster word list" for strengthening/weakening the effect of the following word
- interpreting repeated letters ("baaaaaad"), CAPITALS and !!!
- idioms
- negation
- ...

Thelwall, M., Buckley, K., & Paltoglou, G. (2012). Sentiment strength detection for the social Web. *Journal of the American Society for Information Science and Technology*, 63(1), 163-173.

1	abandon*	-2	liwc uness specified otherwise						
2	abate -2		General Inquirer Feb 2010						
3	abdicate*	-2	General Inquirer Feb 2010						
4	abhor* -4		General Inquirer Feb 2010						
5	abject -2		General Inquirer Feb 2010						
6	abnormal*	-2	General Inquirer Feb 2010						
7	abolish*	-2	General Inquirer Feb 2010						
8	abomina*	-3	General Inquirer Feb 2010						
9	abrasive*	-2	General Inquirer Feb 2010						
10	abrupt -2		General Inquirer Feb 2010						
11	abscond*	-2	General Inquirer Feb 2010						
12	absence -2		General Inquirer Feb 2010						
13	absent* -2		General Inquirer Feb 2010						
14	absurd* -2		Feb-11						
15	abuse* -4								
16	abusi* -4	removed	accept 1 accepta*	2	accepted	2	accepting	2	accepts 2
17	abyss -2		General Inquirer Feb 2010						
18	accident*	-2	General Inquirer Feb 2010						
19	accomplish*	2	Hannes GI add						
20	accost* -2		General Inquirer Feb 2010						
21	accursed	-2	General Inquirer Feb 2010						
22	accus* -2		General Inquirer Feb 2010						
23	accusation*	-2	General Inquirer Feb 2010						
24	ache* -2								
25	achen* 1	kev							
26	acher* 1	kev							
27	acheson 1	kev							
28	acheta 1	kev							
29	aching -2	removed	active* 2						
30	acrimon*	-2	General Inquirer Feb 2010						
31	addict* -2		General Inquirer Feb 2010						

25 lines (25 sloc) | 191 Bytes

```
1  aren't
2  arent
3  can't
4  cannot
5  cant
6  couldn't
7  couldnt
8  didn't
9  didnt
10 doesn't
11 doesnt
12 don't
13 dont
14 hasn't
15 hasnt
16 isn't
17 isnt
18 never
19 not
20 shouldn't
21 shouldnt
22 won't
23 wont
24 wouldn't
25 wouldnt
```

# Advanced approaches

## Take the structure of a text into account

- Try to apply linguistics concepts to identify sentence structure
- can identify negations
- can interpret intensifiers

# Example

```

1 from pattern.nl import sentiment
2 >>> sentiment("Great service by @NSHighspeed")
3 (0.8, 0.75)
4 >>> sentiment("Really")
5 (0.0, 1.0)
6 >>> sentiment("Really Great service by @NSHighspeed")
7 (1.0, 1.0)

```

(polarity, subjectivity) with

$-1 \leq \text{polarity} \leq +1$

$0 \leq \text{subjectivity} \leq +1$  )

Unlike in pure bag-of-words approaches, here, the overall sentiment is not just the sum or the average of its parts!

De Smedt, T., & Daelemans W. (2012). Pattern for Python. *Journal of Machine Learning Research*, 13, 2063-2067.

# Advanced approaches

## pro

- understand intensifiers or negation
- thus: higher accuracy

## con

- Black box? Or do we understand the algorithm?
- Difficult to adapt to own needs
- *really* much better results?



# Data analysis 1: Sentiment analysis

## A sentiment analysis tailored to your needs!

# A sentiment analysis tailored to your needs!

## Identifying suicidal texts

- Bag-of-words-approach with very specific dictionary
- added negation
- added regular expression search for key phrases
- Very specific design requirements: False positives are OK, false negatives not!

Huang, Y.-P., Goh, T., & Liew, C.L. (2007). Hunting suicide notes in web 2.0 – preliminary findings. *Ninth IEEE International Symposium on Multimedia*. Retrieved from <http://ieeexplore.ieee.org/stamp/stamp.jsp?tp=&arnumber=4476021>

group suicide notes by these characteristics:

- Pestian, J.P.; Matykiewicz, P., Linn-Gust, M., South, B., Uzuner, O., Wiebe, J., Cohen, K.B., Hurdle, J., & Brew, C. (2012). Sentiment analysis of suicide notes: A shared task. *Biomedical Informatics Insights*, 5(1), p. 3-16. Retrieved from <http://europepmc.org/articles/PMC3200408?pdf=render>

Packages for sentiment analysis

# Which packages are easy to use?

vader pro: in NLTK module, con: English only

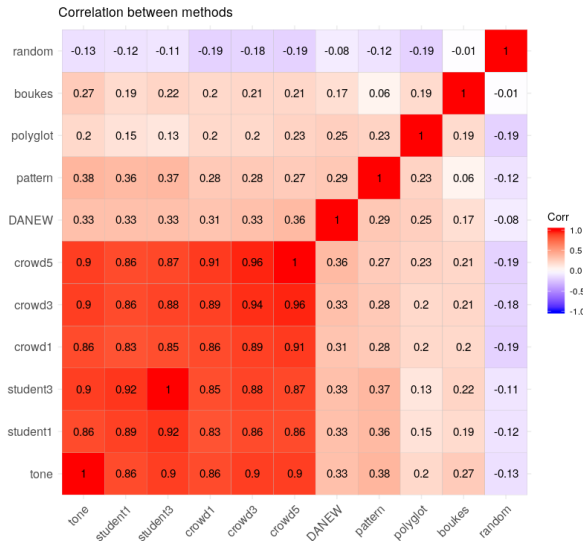
pattern pro: multiple languages (including Dutch)

sentistrength pro: multiple languages, widely used, con: needs  
Python wrapper, license

vader: Chapter 6.3; pattern: Chapter 6.5; sentistrength: Chapter 6.4

**BUT: Keep in mind that the results of *any*  
off-the-shelf-package might be biased and/or noisy in *your*  
domain!**

## Packages for sentiment analysis



Note: student3 (and crowd3, crowd5) are the majority vote between 3 (or 5) student/crowd coders.  
student1, crowd1, and crowd3 are summary values for multiple (combinations of) coders,  
so the diagonal reflects the average correlation between them

Boukes, M., van der Velde, R.N., & Vliegthart, R. (2018). The good and bad in economic news: Comparing (automatic) measurements of sentiment in Dutch economic news. *International Communication Association*

# A possible recipe for doing your sentiment analysis

- ❶ Construct a list `data` of strings with your input data
- ❷ Create an empty list `sent` for storing the results
- ❸ For each text `t` in `data`, estimate the sentiment of `t` and append the result to `sent`<sup>1</sup>
- ❹ Confirm that `len(data) == len(sent)`
- ❺ use `zip()` and a `csv.writer` to write input and output next to each other to a csv file.

---

<sup>1</sup>use multiple lists instead if you estimate for instance subjectivity *and* polarity

# Supervised ML ( $\Rightarrow$ week 7)

An alternative state-of-the-art approach:

## Use supervised machine learning

- Instead of defining rules, hand-code (“annotate”) the sentiment of some tweets manually and let the computer find out which words or characters (“features”) predict sentiment
- Then use this model to predict sentiment for other tweets
- Essentially the same like what you know since the second year of your Bachelor: regression analysis (but now with DV sentiment and IV’s word occurrences)

Gonzalez-Bailon, S., & Paltoglou, G. (2015). Signals of public opinion in online communication: A comparison of methods and data sources. *The ANNALS of the American Academy of Political and Social Science*, 659(1), 95–107.



- Take-home message
- Mid-term take-home exam
- Next meetings

## Take-home messages

## What you should be familiar with:

- You should have *completely* understood last week's exercise. Re-read it if necessary.
- Approaches to the analysis (e.g., structure vs. content)
- Types of sentiment analysis, application areas, pros and cons



