# Working with text: Topic modelling

## Damian Trilling

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

Afdeling Communicatiewetenschap
Universiteit van Amsterdam

18–1–2018

## Today

**1** Automated Content Analysis (ACA)

**2** Basic top-down ACA: Dictionary- and string-based methods
   Regular expressions

**3** Unsupervised Machine Learning
   PCA
   LDA

**4** Example and exercise

What's Automated Content Analysis?

**Methodological approach**

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

deductive                                                                 inductive

Boumans, J. W., & Trilling, D. (2016). Taking stock of the toolkit: An overview of relevant autmated content analysis approaches and techniques for digital journalism scholars. *Digital Journalism, 4*(1), 8–23. doi:10.1080/21670811.2015.1096598

**Basic ACA: Dictionary- and string-based methods**
Regular expressions

# Regular Expressions: What and why?

## What is a regexp?

- a *very* widespread way to describe patterns in strings

# 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

Automated Content Analysis  **Basic ACA**  Unsupervised Machine Learning  Supervised Machine Learning  Example and exercise
○●○○○○○○       ○○○○○○○○○○

Regular expressions

# 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 text editors (!), in STATA, R, Python, . . .

# An example

We wanted to find references to companies in several years of news coverage

Problems:

- Spelling variations (ABN, ABN Amro, ABN-Amro, . . . )
- Shouldn't be in the middle of the word, but *can* be at the beginning of a word, optionally connected with a hyphen ("ABN-topman", "Shellstation")

For instance,
\bING(?:-.*?)?\b
allows to specify exactly this.

Strycharz, J., Strauss, N., & Trilling, D. (2017). The role of media coverage in explaining stock market fluctuations: insights for strategic financial communication. *International Journal of Strategic Communication, online first.* doi:10.1080/1553118X.2017.1378220

Jonkman, J. G., Trilling, D., Verhoeven, P., & Vliegenthart, R. (2016). More or less diverse: An assessment of the effect of attention to media salient company types on media agenda diversity in Dutch newspaper coverage between 2007 and 2013. *Journalism, online first.* doi:10.1177/1464884916680371

# Basic regexp elements

### Alternatives

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

Twitter|Facebook matches either Twitter or Facebook

. matches any character

# 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

# 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 (!)

This is **top-down**: we determined a priori what to look for. But what if we do not want to make such assumptions but want to

look what 'emerges' from the data? Enter **bottom-up** approaches:

Unsupervised Machine Learning

**Methodological approach**

|  | Counting and Dictionary | Supervised Machine Learning | Unsupervised Machine Learning |
|---|---|---|---|
| **Typical research interests and content features** | visibility analysis sentiment analysis subjectivity analysis | frames topics gender bias | frames topics |
| **Common statistical procedures** | string comparisons counting | support vector machines naive 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 autmated content analysis approaches and techniques for digital journalism scholars. *Digital Journalism, 4*(1), 8–23. doi:10.1080/21670811.2015.1096598

# Supervised vs. unsupervised learning

## Unsupervised

- No manually coded data

- We want to identify patterns or to make groups of most similar cases

# Supervised vs. unsupervised learning

### Unsupervised

- No manually coded data
- We want to identify patterns or to make groups of most similar cases

Example: We have a dataset of Facebook-massages on an organizations' page. We use clustering to group them and later interpret these clusters (e.g., as complaints, questions, praise, . . . )

# Supervised vs. unsupervised learning

## Unsupervised

- No manually coded data
- We want to identify patterns or to make groups of most similar cases

Example: We have a dataset of Facebook-massages on an organizations' page. We use clustering to group them and later interpret these clusters (e.g., as complaints, questions, praise, . . . )

## Supervised

- We code a small dataset by hand and use it to "train" a machine
- The machine codes the rest

# Supervised vs. unsupervised learning

## Unsupervised

- No manually coded data
- We want to identify patterns or to make groups of most similar cases

Example: We have a dataset of Facebook-massages on an organizations' page. We use clustering to group them and later interpret these clusters (e.g., as complaints, questions, praise, . . . )

## Supervised

- We code a small dataset by hand and use it to "train" a machine
- The machine codes the rest

Example: We have 2,000 of these messages grouped into such categories by human coders. We then use this data to group all remaining messages as well.

inductive and bottom-up:
**unsupervised machine learning**

inductive and bottom-up:
**unsupervised machine learning**

(something you aready did in your Bachelor – no kidding.)

Automated Content Analysis   Basic ACA   **Unsupervised Machine Learning**   Supervised Machine Learning   Example and exercise
0000000   ●000000000

PCA

Principal Component Analysis? How does *that* fit in here?

# Principal Component Analysis? How does *that* fit in here?

In fact, PCA is used everywhere, even in image compression

# Principal Component Analysis? How does *that* fit in here?

## PCA in ACA

- Find out what word cooccur (inductive frame analysis)

- Basically, transform each document in a vector of word frequencies and do a PCA

Automated Content Analysis   Basic ACA   **Unsupervised Machine Learning**   Supervised Machine Learning   Example and exercise
ooooooo   oooooooooo

PCA

# A so-called term-document-matrix

```
1          w1,w2,w3,w4,w5,w6 ...
2    text1, 2, 0, 0, 1, 2, 3 ...
3    text2, 0, 0, 1, 2, 3, 4 ...
4    text3, 9, 0, 1, 1, 0, 0 ...
5    ...
```

Automated Content Analysis  Basic ACA  **Unsupervised Machine Learning**  Supervised Machine Learning  Example and exercise
          0000000      0●00000000

PCA

# A so-called term-document-matrix

```
1        w1,w2,w3,w4,w5,w6 ...
2   text1, 2, 0, 0, 1, 2, 3 ...
3   text2, 0, 0, 1, 2, 3, 4 ...
4   text3, 9, 0, 1, 1, 0, 0 ...
5   ...
```

These can be simple counts, but also more advanced metrics, like
tf-idf scores (where you weigh the frequency by the number of
documents in which it occurs), cosine distances, etc.

Automated Content Analysis  Basic ACA  **Unsupervised Machine Learning**  Supervised Machine Learning  Example and exercise
                0000000    000●000000

PCA

# PCA: implications and problems

- given a term-document matrix, easy to do with any tool
- probably extremely skewed distributions
- some problematic assumptions: does the goal of PCA, to find a solution in which one word loads on *one* component match real life, where a word can belong to several topics or frames?

Automated Content Analysis  Basic ACA  **Unsupervised Machine Learning**  Supervised Machine Learning  Example and exercise
0000000  000●000000

LDA

Enter **topic modeling with Latent Dirichlet Allocation (LDA)**

# LDA, what's that?

## No mathematical details here, but the general idea

- There are $k$ topics, $T_1 \ldots T_k$
- Each document $D_i$ consists of a mixture of these topics, e.g. $80\% T_1, 15\% T_2, 0\% T_3, \ldots 5\% T_k$
- On the next level, each topic consists of a specific probability distribution of words
- Thus, based on the frequencies of words in $D_i$, one can infer its distribution of topics
- Note that LDA (like PCA) is a Bag-of-Words (BOW) approach

# Doing a LDA in Python

You can use gensim (Řehůřek & Sojka, 2010) for this.

Řehůřek, R., & Sojka, P. (2010). Software framework for topic modelling with large corpora. *Proceedings of the LREC 2010 Workshop on New Challenges for NLP Frameworks*, pp. 45–50. Valletta, Malta: ELRA.

```
 1    from gensim import corpora, models
 2
 3    NTOPICS = 100
 4    LDAOUTPUTFILE="topicscores.tsv"
 5
 6    # Create a BOW represenation of the texts
 7    id2word = corpora.Dictionary(texts)
 8    mm =[id2word.doc2bow(text) for text in texts]
 9
10    # Train the LDA models.
11    lda = models.ldamodel.LdaModel(corpus=mm, id2word=id2word, num_topics=
         NTOPICS)
12
13    # Print the topics.
14    for top in lda.print_topics(num_topics=NTOPICS, num_words=5):
15        print ("\n",top)
16
17    # save topic scores
18    scoresperdoc=lda.inference(mm)
19    with open(LDAOUTPUTFILE,"w",encoding="utf-8") as fo:
20      for row in scoresperdoc[0]:
21        fo.write("\t".join(["{:0.3f}".format(score) for score in row]))
22        fo.write("\n")
```

Automated Content Analysis  Basic ACA  **Unsupervised Machine Learning**  Supervised Machine Learning  Example and exercise
0000000  0000000●00

LDA

## Output: Topics (below) & topic scores (next slide)

```
 1  0.069*fusie + 0.058*brussel + 0.045*europesecommissie + 0.036*europese +
        0.023*overname
 2  0.109*bank + 0.066*britse + 0.041*regering + 0.035*financien + 0.033*
        minister
 3  0.114*nederlandse + 0.106*nederland + 0.070*bedrijven + 0.042*rusland +
        0.038*russische
 4  0.093*nederlandsespoorwegen + 0.074*den + 0.036*jaar + 0.029*onderzoek +
        0.027*raad
 5  0.099*banen + 0.045*jaar + 0.045*productie + 0.036*ton + 0.029*aantal
 6  0.041*grote + 0.038*bedrijven + 0.027*ondernemers + 0.023*goed + 0.015*
        jaar
 7  0.108*werknemers + 0.037*jongeren + 0.035*werkgevers + 0.029*jaar +
        0.025*werk
 8  0.171*bank + 0.122* + 0.041*klanten + 0.035*verzekeraar + 0.028*euro
 9  0.162*banken + 0.055*bank + 0.039*centrale + 0.027*leningen + 0.024*
        financiele
10  0.052*post + 0.042*media + 0.038*nieuwe + 0.034*netwerk + 0.025*
        personeel
11  ...
```

Edit | Browse

Filter | Variables | Properties | Snapshots

topic4[2] | .019

| | source2 | firstwords | polarity | subjectivity | pubdate_day | pubdate_mo~h | pubdate_year | pubdate_da~k | topic1 | topic2 | topic3 | topic4 | topic5 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | nrc handelsblad | palingsound schinke | -.0086207 | .6069971 | 31 | 12 | 2011 | zaterdag | .018 | .019 | 3.587 | .019 | .019 |
| 2 | nrc handelsblad | groep investeerders | -.1041667 | .3129167 | 31 | 12 | 2011 | zaterdag | .018 | .019 | .019 | .019 | .019 |
| 3 | nrc handelsblad | abnamro debacle ij | .0082292 | .4895443 | 31 | 12 | 2011 | zaterdag | .018 | 27.71 | .019 | .019 | .019 |
| 4 | nrc handelsblad | abnamro financi`l~ | -.0179617 | .5706419 | 31 | 12 | 2011 | zaterdag | 15.1 | .019 | 2.646 | .019 | .019 |
| 5 | nrc handelsblad | crisis verhouding k | .0758049 | .5348864 | 31 | 12 | 2011 | zaterdag | .018 | .019 | 9.008 | .019 | .019 |
| 6 | nrc handelsblad | snel vakantie vrije | -.016315 | .5118800 | 31 | 12 | 2011 | zaterdag | .018 | .019 | .019 | .019 | .019 |
| 7 | nrc handelsblad | herinnering dees le | .18075 | .6203333 | 31 | 12 | 2011 | zaterdag | .018 | .019 | .019 | .019 | .019 |
| 8 | nrc handelsblad | hackers publiceren | .1454545 | .4545455 | 31 | 12 | 2011 | zaterdag | .018 | .019 | .019 | .019 | .019 |
| 9 | nrc handelsblad | waterballet montevi | -.2333333 | .4333333 | 31 | 12 | 2011 | zaterdag | .018 | .019 | .019 | .019 | .019 |
| 10 | nrc handelsblad | bouw dupe ambities | .0925417 | .5939167 | 5 | 11 | 2010 | vrijdag | .018 | .019 | .078 | 2.442 | .019 |
| 11 | nrc handelsblad | eindelijk wint nuch | .1755893 | .48125 | 5 | 11 | 2010 | vrijdag | .018 | .019 | 8.302 | .019 | .019 |
| 12 | nrc handelsblad | oud nieuws tv bbcst | .02 | .4322222 | 5 | 11 | 2010 | vrijdag | .018 | 10.853 | .019 | .019 | .019 |
| 13 | nrc handelsblad | tmg hyves krantenbe | .0425283 | .5420412 | 5 | 11 | 2010 | vrijdag | .018 | .019 | .019 | .019 | .019 |
| 14 | nrc handelsblad | getuigenis rechter | .0858929 | .5770833 | 5 | 11 | 2010 | vrijdag | .018 | .019 | .019 | 11.621 | .019 |
| 15 | nrc handelsblad | akzonobel philips g | .0220455 | .4381818 | 5 | 11 | 2010 | vrijdag | .018 | .019 | .019 | .019 | .019 |
| 16 | nrc handelsblad | mondiaal kritiek be | -.038172 | .3094624 | 5 | 11 | 2010 | vrijdag | .018 | 19.957 | .019 | .019 | .019 |
| 17 | nrc handelsblad | export diamant fiat | .0628571 | .4438095 | 5 | 11 | 2010 | vrijdag | .018 | 4.745 | .019 | .019 | .019 |
| 18 | nrc handelsblad | canada bod potash r | .0252924 | .4795322 | 5 | 11 | 2010 | vrijdag | .018 | 26.741 | .019 | .019 | .019 |
| 19 | nrc handelsblad | zwakke bouwsector c | -.0171 | .4736333 | 14 | 3 | 2009 | NA | .018 | .019 | .019 | .019 | 4.806 |
| 20 | nrc handelsblad | pensioenconflict wa | .028114 | .4636842 | 14 | 3 | 2009 | NA | .018 | .019 | .019 | .019 | .019 |
| 21 | nrc handelsblad | rechter allim loon | .1318182 | .3939394 | 14 | 3 | 2009 | NA | .018 | .019 | .019 | .019 | .019 |
| 22 | nrc handelsblad | bad bank remedie da | .0891026 | .550641 | 14 | 3 | 2009 | NA | .018 | 10.235 | .019 | .019 | .019 |
| 23 | nrc handelsblad | bescheiden salaris | -.075 | .56 | 14 | 3 | 2009 | NA | .018 | .019 | .019 | .019 | .019 |
| 24 | nrc handelsblad | generalmotors autos | .0138889 | .4388889 | 14 | 3 | 2009 | NA | .018 | .019 | .019 | .019 | .019 |
| 25 | nrc handelsblad | rusland rozen tuinb | .0314141 | .5643951 | 14 | 3 | 2009 | NA | .018 | .019 | 24.595 | .019 | .019 |
| 26 | nrc handelsblad | cynisme oplossing k | .0100833 | .6511667 | 14 | 3 | 2009 | NA | .018 | .019 | .019 | .019 | .019 |
| 27 | nrc handelsblad | the good bed ugly l | .0265504 | .5298449 | 13 | 3 | 2009 | NA | .018 | .019 | .019 | .019 | .019 |
| 28 | nrc handelsblad | kerk stroom nietswe | -.0087719 | .6149123 | 13 | 3 | 2009 | NA | .018 | .019 | .019 | .019 | .019 |
| 29 | nrc handelsblad | kerk stroom gaud ac | 0 | 0 | 13 | 3 | 2009 | NA | .018 | .019 | .019 | .019 | .019 |
| 30 | nrc handelsblad | supersnelle koekenp | 0 | 0 | 13 | 3 | 2009 | NA | .018 | .019 | .019 | .019 | .019 |
| 31 | nrc handelsblad | dalailama chinese e | 0 | 0 | 13 | 3 | 2009 | NA | .018 | .019 | .019 | .019 | .019 |
| 32 | nrc handelsblad | bezuinigen hulpgeld | .0894192 | .4560606 | 4 | 10 | 2009 | NA | .018 | .019 | .019 | .019 | .019 |
| 33 | nrc handelsblad | vaders arbeidsethos | .0160985 | .5575758 | 4 | 10 | 2009 | NA | .018 | .019 | .019 | .019 | .019 |
| 34 | nrc handelsblad | varkens lux winnaar | .040873 | .6218254 | 4 | 10 | 2008 | NA | .018 | .019 | .019 | .019 | .019 |
| 35 | nrc handelsblad | liberale kinderopva | .1179095 | .5297855 | 4 | 10 | 2008 | NA | .018 | .019 | .019 | .019 | 1.03 |
| 36 | nrc handelsblad | banken verzinsels k | .068521 | .6300389 | 4 | 10 | 2008 | NA | 8.232 | .019 | .019 | .019 | .019 |
| 37 | nrc handelsblad | rabobanktopman bert | 0 | 0 | 4 | 10 | 2008 | NA | .018 | .019 | .019 | .019 | .019 |
| 38 | nrc handelsblad | kinderopvang bril v | 0 | 0 | 4 | 10 | 2008 | NA | .018 | .019 | .019 | .019 | .019 |
| 39 | nrc handelsblad | tassen gevoel verli | 0 | 0 | 4 | 10 | 2008 | NA | .018 | .019 | .019 | .019 | .019 |
| 40 | nrc handelsblad | abnamro winkelend p | .0876761 | .62277 | 4 | 10 | 2008 | NA | .018 | .019 | 6.904 | .019 | 5.511 |
| 41 | nrc handelsblad | abnamro belgiv` mole | .0439506 | .4976852 | 4 | 10 | 2008 | NA | .018 | .019 | .019 | .019 | .019 |
| 42 | nrc handelsblad | abnamro handen deut | .1838401 | .5264302 | 4 | 10 | 2008 | NA | .018 | .019 | 1.854 | .019 | .019 |
| 43 | nrc handelsblad | abnamro fortis bank | .0842391 | .494058 | 4 | 10 | 2008 | NA | 4.939 | .019 | 14.39 | .019 | .019 |
| 44 | nrc handelsblad | abnamro fortis spra | .0540715 | .6290807 | 4 | 10 | 2008 | NA | .018 | .019 | .019 | .019 | .019 |
| 45 | nrc handelsblad | abnamro fortis jaar | .0297297 | .4960135 | 4 | 10 | 2008 | NA | .018 | 11.041 | .019 | .019 | .019 |
| 46 | nrc handelsblad | abnamro nederland s | .1006944 | .6450695 | 4 | 10 | 2008 | NA | .018 | .019 | .019 | .019 | .019 |
| 47 | nrc handelsblad | abnamro belgiv` mole | .0405952 | .5004464 | 4 | 10 | 2008 | NA | 6.730535 | .019 | .019 | .019 | 12.682 |
| 48 | nrc handelsblad | arbeidsmarkt vs sle | .0166667 | .4 | 4 | 10 | 2008 | NA | 7.103 | .019 | .019 | .019 | 12.682 |

Variables

Q- Enter filter text here

| | Name | Label |
|---|---|---|
| | byline | |
| | section | |

Properties

▼ Variables
Name | section
Label |
Type | str26
Format | %26s
Value Label |
Notes |
▼ Data
▶ Filename | topicscores.
Label |
Notes |
Variables | 164
Observations | 28,406
Size | 14.06M
Memory | 64M

Vars: 158 of 164 | Obs: 28,406 | | Filter: Off

Automated Content Analysis  Basic ACA  **Unsupervised Machine Learning**  Supervised Machine Learning  Example and exercise
○○○○○○○  ○○○○○○○○●

LDA

# Visualization with pyldavis

```
1  import pyLDAvis
2  import pyLDAvis.gensim
3  % first estiate gensim model, then:
4  vis_data = pyLDAvis.gensim.prepare(lda,mm,id2word)
5  pyLDAvis.display(vis_data)
```

**Supervised machine learning** is something for another time . . .

Example and exercise

Let's have a look at the EU-speech dataset (Jupyter Notebook).
I'll first walk you through the example, afterwards, you have time
to play with the data yourself.

# Damian Trilling

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