# Unsupervised approaches to detecting topics in news and politics

Damian Trilling

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

Afdeling Communicatiewetenschap
Universiteit van Amsterdam

18 October 2019, Salamanca

# Today

**1** Types of Automated Content Analysis

**2** PCA and Clustering

**3** LDA Topic models
     An introduction to LDA
     Choosing the best (or a good) topic model
     Using topic models
     Other forms of topic models

**4** Let's do it!

Types of 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**

# Some terminology

## Supervised machine learning

You have a dataset with both predictor and outcome (independent and dependent variables; features and labels) — a *labeled* dataset.

# Some terminology

## Supervised machine learning

You have a dataset with both predictor and outcome (independent and dependent variables; features and labels) — a *labeled* dataset.   Think of regression: You measured x1, x2, x3 and you want to predict y, which you also measured

# Some terminology

## Supervised machine learning

You have a dataset with both predictor and outcome (independent and dependent variables; features and labels) — a *labeled* dataset.

## Unsupervised machine learning

You have no labels.

Types of Automated Content Analysis
○○●

PCA and Clustering
○○○○○

LDA Topic models
○○○○○○○○○○○○○○○○○○

Let's do it!
○

Top-down vs. bottom-up

# Some terminology

## Supervised machine learning

You have a dataset with both predictor and outcome (independent and dependent variables; features and labels) — a *labeled* dataset.

## Unsupervised machine learning

You have no labels. (You did not measure y)

Types of Automated Content Analysis         PCA and Clustering         LDA Topic models         Let's do it!

Top-down vs. bottom-up

# Some terminology

## Unsupervised machine learning

You have no labels.

**You probably already know some techniques to find out how `x1`, `x2`,... `x_i` co-occur:**

- Principal Component Analysis (PCA) and Singular Value Decomposition (SVD)

- Cluster analysis

- ...

PCA and Clustering

| Types of Automated Content Analysis | PCA and Clustering | LDA Topic models | Let's do it! |
|---|---|---|---|
| ooo | o●ooo | oooooooooooooooo | o |

PCA and Clustering

Let's assume we want to find out the topics in a large corpus of documents

We could, for example

- use PCA to find out related features (and interpret those as topics)
- or use clustering to find similar documents (and then look at the words they share to interpret as topics)

Let's assume we want to find out the topics in a large corpus of documents

We could, for example

- use PCA to find out related features (and interpret those as topics)
- or use clustering to find similar documents (and then look at the words they share to interpret as topics)

Actually, we have *two* things we want to model:

❶ Which topics can we extract from the corpus?

❷ How present is each of these topics in each text in the corpus?

| Types of Automated Content Analysis | PCA and Clustering | LDA Topic models | Let's do it! |
|---|---|---|---|
| ○○○ | ○○●○○ | ○○○○○○○○○○○○○○○ | ○ |

PCA and Clustering

# Recap: PCA

Document-term 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.

# Recap: PCA

### Document-term 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.

- 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?

# Recap: clustering

- given a term-document matrix, we can easily find clusters of documents that resemble each other
- but also here does the goal of cluster analysis, assigning each document to *one* cluster, match real life?

# We need other models to

**1** model *simultaneously* (a) which topics we find in the whole corpus, and (b) which of these topics are present in which document; while at the same time

Types of Automated Content Analysis | PCA and Clustering | LDA Topic models | Let's do it!
000 | 00000 | 0000000000000000 | 0

PCA and Clustering

# We need other models to

1. model *simultaneously* (a) which topics we find in the whole corpus, and (b) which of these topics are present in which document; while at the same time

2. allowing (a) words to be part of multiple topics, and (b) multiple topics to be present in one document; and

Maier, D., Waldherr, A., Miltner, P., Wiedemann, G., Niekler, A., Keinert, A., . . . Adam, S. (2018). Applying LDA Topic Modeling in Communication Research: Toward a Valid and Reliable Methodology. *Communication Methods and Measures, 12*(2–3), 93–118. doi:10.1080/19312458.2018.1430754

Types of Automated Content Analysis | PCA and Clustering | LDA Topic models | Let's do it!
000 | 00000 | 0000000000000000 | 0

PCA and Clustering

# We need other models to

❶ model *simultaneously* (a) which topics we find in the whole corpus, and (b) which of these topics are present in which document; while at the same time

❷ allowing (a) words to be part of multiple topics, and (b) multiple topics to be present in one document; and

❸ being able to make connections between words "even if they never actually occured in a document together" (Maier et al, 2018, p. 96)

Maier, D., Waldherr, A., Miltner, P., Wiedemann, G., Niekler, A., Keinert, A., ... Adam, S. (2018). Applying LDA Topic Modeling in Communication Research: Toward a Valid and Reliable Methodology. *Communication Methods and Measures, 12*(2–3), 93–118. doi:10.1080/19312458.2018.1430754

Types of Automated Content Analysis        PCA and Clustering        LDA Topic models        Let's do it!
000                                         00000                     ●000000000000000        0
An introduction to LDA

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

Types of Automated Content Analysis   PCA and Clustering   LDA Topic models   Let's do it!
000                                    00000                 0●00000000000000    0
An introduction to 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

Types of Automated Content Analysis              PCA and Clustering              **LDA Topic models**              Let's do it!

000                              00000                            000●0000000000000             O

An introduction to LDA

# Doing a LDA in Python

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

Let us assume you have a list of lists of words (!) called `texts`:

```
1  articles=['The tax deficit is higher than expected. This said xxx ...',
       'Germany won the World Cup. After a']
2  texts=[art.split() for art in articles]
```

which looks like this:

```
1  [['The', 'tax', 'deficit', 'is', 'higher', 'than', 'expected.', 'This',
       'said', 'xxx', '...'], ['Germany', 'won', 'the', 'World', 'Cup.', '
       After', 'a']]
```

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.

```python
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  mylda = models.ldamodel.LdaModel(corpus=mm, id2word=id2word, num_topics=
        NTOPICS, alpha="auto")
12
13  # Print the topics.
14  for top in mylda.print_topics(num_topics=NTOPICS, num_words=5):
15    print ("\n",top)
16
17  print ("\nFor further analysis, a dataset with the topic score for each
        document is saved to",LDAOUTPUTFILE)
18
19  scoresperdoc=mylda.inference(mm)
20
21  with open(LDAOUTPUTFILE,"w",encoding="utf-8") as fo:
22    for row in scoresperdoc[0]:
23      fo.write("\t".join(["{:0.3f}".format(score) for score in row]))
24  fo.write("\n")
```

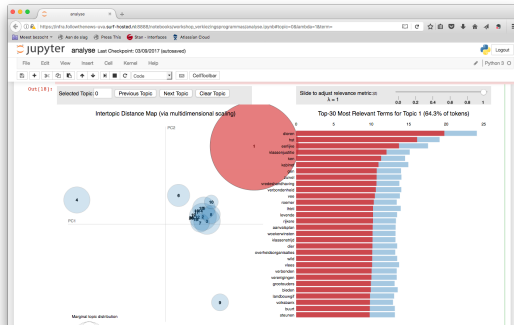| Types of Automated Content Analysis | PCA and Clustering | LDA Topic models | Let's do it! |
|---|---|---|---|
| ○○○ | ○○○○○ | ○○○○●○○○○○○○○○○○ | ○ |

An introduction to 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 | .4095443 | 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 | -.016135 | .5118800 | 31 | 12 | 2011 | zaterdag | .018 | .019 | .019 | .019 | .019 |
| 7 | nrc handelsblad | herinnering doos le | .18875 | .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 | .1755093 | .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 hvves krantenbe | .0425203 | .5420412 | 5 | 11 | 2010 | vrijdag | .018 | .019 | .019 | .019 | .019 |
| 14 | nrc handelsblad | getuigenis rechter | .0858929 | .5770833 | 5 | 11 | 2010 | vrijdag | .018 | .019 | 11.621 | .019 | .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 l | .0252924 | .4795322 | 5 | 11 | 2010 | vrijdag | .018 | 26.741 | .019 | .019 | .019 |
| 19 | nrc handelsblad | zwakke bouwsector c | -.0171 | .4733333 | 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 | .5643051 | 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 | 2008 | NA | .018 | .019 | .019 | .019 | .019 |
| 33 | nrc handelsblad | vaders arbeidsethos | .0160985 | .5575758 | 4 | 10 | 2008 | 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 | 1.03 | .019 |
| 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 werli | 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 | .6035455 | 4 | 10 | 2008 | NA | .018 | .019 | .019 | .019 | .019 |
| 47 | nrc handelsblad | abnamro belgiv' mole | .0405952 | .5804464 | 4 | 10 | 2008 | NA | .018 | .019 | .019 | .019 | .019 |
| 48 | nrc handelsblad | arbeidsmarkt vs sle | .0166667 | .4 | 4 | 10 | 2008 | NA | 7.103 | .019 | .019 | .019 | 12.682 |

Variables

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

Types of Automated Content Analysis    PCA and Clustering    LDA Topic models    Let's do it!
○○○                                      ○○○○○               ○○○○○○○●○○○○○○○○        ○
An introduction to LDA

# Visualization with pyldavis

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

Types of Automated Content Analysis    PCA and Clustering    LDA Topic models    Let's do it!
000                                     00000                 0000000●00000000    0
An introduction to LDA

# Visualization with pyldavis

Short note about the $\lambda$ setting:
It influences the ordering of the words in pyldavis.

> "For $\lambda = 1$, the ordering of the top words is equal to the
> ordering of the standard conditional word probabilities. For
> $\lambda$ close to zero, the most specific words of the topic will
> lead the list of top words. In their case study, Sievert and
> Shirley (2014, p. 67) found the best interpretability of
> topics using a $\lambda$-value close to .6, which we adopted for
> our own case" (Maier et al., 2018, p. 107)

Maier, D., Waldherr, A., Miltner, P., Wiedemann, G., Niekler, A., Keinert, A., . . . Adam, S. (2018). Applying LDA
Topic Modeling in Communication Research: Toward a Valid and Reliable Methodology. *Communication Methods
and Measures, 12*(2–3), 93–118. doi:10.1080/19312458.2018.1430754

Types of Automated Content Analysis     PCA and Clustering     LDA Topic models     Let's do it!
○○○                                      ○○○○○                  ○○○○○○○○○●○○○○○○○     ○
Choosing the best (or a good) topic model

# Choosing the best (or a good) topic model

- There is no single best solution (e.g., do you want more coarse or fine-grained topics?)
- Non-deterministic
- Very sensitive to preprocessing choices
- Interplay of both metrics and (qualitative) interpretability

See for more elaborate guidance:
Maier, D., Waldherr, A., Miltner, P., Wiedemann, G., Niekler, A., Keinert, A., ... Adam, S. (2018). Applying LDA Topic Modeling in Communication Research: Toward a Valid and Reliable Methodology. *Communication Methods and Measures, 12*(2–3), 93–118. doi:10.1080/19312458.2018.1430754

Types of Automated Content Analysis        PCA and Clustering        LDA Topic models        Let's do it!
000                                          00000                    000000000000000000    0
Choosing the best (or a good) topic model

# Evaluation metrics (closer to zero is better)

### perplexity

A goodness-of-fit measure, answering the question: If we do a train-test split, how well does the trained model fit the test data?

# Evaluation metrics (closer to zero is better)

### perplexity

A goodness-of-fit measure, answering the question: If we do a train-test split, how well does the trained model fit the test data?

### coherence

- mean coherence of the whole model: attempts to quantify the interpretability
- coherence per topic: allows to get topics that are most likely to be coherently interpreted (`.top_topics()`)

Types of Automated Content Analysis        PCA and Clustering        LDA Topic models        Let's do it!
000                                          00000                     000000000000000000000   0
Choosing the best (or a good) topic model

# Choosing $k$: How many topics do we want?

- Typical values: $10 < k < 200$
- Too low: losing nuance, so broad it becomes meaningless
- Too high: picks up tiny pecularities instead of finding general patterns
- There is no inherent ordering of topics (unlike PCA!)
- We can throw away or merge topics later, so if out of $k = 50$ topics 5 are not interpretable and a couple of others overlap, it still may be a good model

# Choosing $\alpha$: how sparse should the document-topic distribution $\theta$ be?

- The higher $\alpha$, the more topics per document
- Default: $1/k$
- But: We can explicitly change it, or – really cool – even learn $\alpha$ from the data (alpha = "auto")

| Types of Automated Content Analysis | PCA and Clustering | LDA Topic models | Let's do it! |
|---|---|---|---|
| ○○○ | ○○○○○ | ○○○○○○○○○○○○○●○○○○ | ○ |

Choosing the best (or a good) topic model

# Choosing $\alpha$: how sparse should the document-topic distribution $\theta$ be?

- The higher $\alpha$, the more topics per document
- Default: $1/k$
- But: We can explicitly change it, or – really cool – even learn $\alpha$ from the data (alpha = "auto")

Takeaway: It takes longer, but you probably want to learn alpha from the data, using multiple passes:

```
1   mylda LdaModel(corpus=tfidfcorpus[ldacorpus], id2word=id2word,
        num_topics=50, alpha='auto', passes=10)
```

| Types of Automated Content Analysis | PCA and Clustering | LDA Topic models | Let's do it! |
| 000 | 00000 | 000000000000000000 | 0 |

Choosing the best (or a good) topic model

# Choosing $\eta$: how sparse should the topic-word distribution $\lambda$ be?

- Can be used to boost specific words
- Can also be learned from the data

Types of Automated Content Analysis    PCA and Clustering    LDA Topic models    Let's do it!
000                                     00000                 0000000000000●000    0
Choosing the best (or a good) topic model

Choosing $\eta$: how sparse should the topic-word distribution $\lambda$ be?

- Can be used to boost specific words
- Can also be learned from the data

Takeaway: Even though you can do eta="auto", this usually does not help you much.

| Types of Automated Content Analysis | PCA and Clustering | LDA Topic models | Let's do it! |
| :--- | :--- | :--- | :--- |
| ○○○ | ○○○○○ | ○○○○○○○○○○○○○●○○ | ○ |

Using topic models

# Using topic models

You got your model – what now?

1. Assign topic scores to documents
2. Label topics
3. Merge topics, throw away boilerplate topics and similar (manually, or aided by cluster analysis)
4. Compare topics between, e.g., outlets
5. or do some time-series analysis.

Example: Tsur, O., Calacci, D., & Lazer, D. (2015). A Frame of Mind: Using Statistical Models for Detection of Framing and Agenda Setting Campaigns. *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing* (pp. 1629–1638).

Other forms of topic models

Types of Automated Content Analysis          PCA and Clustering          LDA Topic models          Let's do it!
000                                          00000                       00000000000000000                o
Other forms of topic models

# Other forms of topic models

- Author-topic models
- Structural topic models
- Non-negative matrix factorization
- . . .

# Let's do it!

(go to Google colab)