

Big Data and Automated Content Analysis (6EC)

Week 6: »Unsupervised machine learning« Monday

Anne Kroon

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

May 8, 2023

UvA RM Communication Science

Today

Unsupervised Machine Learning for Text Classification

- An introduction to LDA

- Choosing the best (or a good) topic model

- Using topic models

Next steps



*Everything clear from previous weeks?
Questions?*

This week, we will get a general overview of working with textual data. Due to a lack of time, I will introduce you to some of the basic concepts, point you to resources, and give you a practical, hands-on introduction.

Boumans and Trilling, 2016: Types of Automated Content Analysis

	Methodological approach		
	<i>Counting and Dictionary</i>	<i>Supervised Machine Learning</i>	<i>Unsupervised Machine Learning</i>
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

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

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 texts = ["Communication in the Digital Society is a very very complex
3 ↪ phenomenon", "I like to study it"]
4 bottom_up = []
5 for t in texts:
6     bottom_up.append(Counter(t.lower().split()).most_common(3))
7     print(bottom_up)
```

This results in:

```
1 [('very', 2), ('Communication', 1), ('in', 1)]
2 [('I', 1), ('like', 1), ('to', 1)]
```

Please note that you can also write this like:

```
1 bottom_up = [Counter(t.split()).most_common(3) for t in texts]
```


A simple bottom-up approach

```
1 from collections import Counter
2 texts = ["Communication in the Digital Society is a very very complex
3 ↪ phenomenon", "I like to study it"]
4 bottom_up = []
5 for t in texts:
6     bottom_up.append(Counter(t.lower().split()).most_common(3))
7     print(bottom_up)
```

This results in:

```
1 [('very', 2), ('Communication', 1), ('in', 1)]
2 [('I', 1), ('like', 1), ('to', 1)]
```

Please note that you can also write this like:

```
1 bottom_up = [Counter(t.split()).most_common(3) for t in texts]
```

A simple bottom-up approach

```
1 from collections import Counter
2 texts = ["Communication in the Digital Society is a very very complex
3 ↪ phenomenon", "I like to study it"]
4 bottom_up = []
5 for t in texts:
6     bottom_up.append(Counter(t.lower().split()).most_common(3))
7     print(bottom_up)
```

This results in:

```
1 [('very', 2), ('Communication', 1), ('in', 1)]
2 [('I', 1), ('like', 1), ('to', 1)]
```

Please note that you can also write this like:

```
1 bottom_up = [Counter(t.split()).most_common(3) for t in texts]
```

A simple top-down approach

```
1 texts = ["Communication in the Digital Society is a very very complex  
↪ phenomenon", "I like to study it"]  
2 features = ["communication", "digital", "study"]  
3 for t in texts:  
4     print(f"\nAnalyzing '{t}':")  
5     for f in features:  
6         print(f"{f} occurs {t.lower().count(f)} times")
```

```
1 Analyzing 'Communication in the Digital Society is a very very complex phenomenon':  
2 communication occurs 1 times  
3 digital occurs 1 times  
4 study occurs 0 times  
5  
6 Analyzing 'I like to study it':  
7 communication occurs 0 times  
8 digital occurs 0 times  
9 study occurs 1 times
```

A simple top-down approach

```
1 texts = ["Communication in the Digital Society is a very very complex  
↪ phenomenon", "I like to study it"]  
2 features = ["communication", "digital", "study"]  
3 for t in texts:  
4     print(f"\nAnalyzing '{t}':")  
5     for f in features:  
6         print(f"{f} occurs {t.lower().count(f)} times")
```

```
1 Analyzing 'Communication in the Digital Society is a very very complex phenomenon':  
2 communication occurs 1 times  
3 digital occurs 1 times  
4 study occurs 0 times  
5  
6 Analyzing 'I like to study it':  
7 communication occurs 0 times  
8 digital occurs 0 times  
9 study occurs 1 times
```



When would you use which approach?

Some considerations

- Both can have a place in your workflow (e.g., bottom-up as first exploratory step)
- You have a clear theoretical expectation? Bottom-up makes little sense.
- But in any case: you need to transform your text into something “countable”.

Some considerations

- Both can have a place in your workflow (e.g., bottom-up as first exploratory step)
- You have a clear theoretical expectation? Bottom-up makes little sense.
- But in any case: you need to transform your text into something “countable”.

Some considerations

- Both can have a place in your workflow (e.g., bottom-up as first exploratory step)
- You have a clear theoretical expectation? Bottom-up makes little sense.
- But in any case: you need to transform your text into something “countable”.

Unsupervised Machine Learning for Text Classification

Supervised vs Unsupervised

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 x_1, x_2, x_3 and you want to predict y , which you also measured

Unsupervised machine learning

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

You might already know some techniques to figure out whether x_1, x_2, \dots, x_i co-occur

- Principal Component Analysis (PCA) and Singular Value Decomposition (SVD)
- Cluster analysis
- Topic modelling (Non-negative matrix factorization and Latent Dirichlet Allocation)

Supervised vs Unsupervised

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 x_1 , x_2 , x_3 and you want to predict y , which you also measured

Unsupervised machine learning

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

You might already know *some* techniques to figure out whether x_1 , x_2, \dots, x_i co-occur

- Principal Component Analysis (PCA) and Singular Value Decomposition (SVD)
- Cluster analysis
- Topic modelling (Non-negative matrix factorization and Latent Dirichlet Allocation)

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 x_1 , x_2 , x_3 and you want to predict y , which you also measured

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

You have a dataset with both predictor and outcome (independent and dependent variables; features and labels) — a *labeled* dataset. Think of regression: You measured x_1 , x_2 , x_3 and you want to predict y , which you also measured

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

Unsupervised Machine Learning for Text Classification

An introduction to 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, $T_1 \dots T_k$
- Each document D_i consists of a mixture of these topics, e.g. $80\% T_1, 15\% T_2, 0\% T_3, \dots 5\% T_k$
- 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 is a Bag-of-Words (BOW) approach

Doing a LDA in Python

We will use gensim (Rehurek10softwareframework) for this (make sure you have version >4.0)

Let us assume you have a list of lists of documents called `texts`:

1

```
print(texts[0][:115])
```

which looks something like:

1

```
'Stop the presses: CNN covered some actual news yesterday when it reported on the story of  
→ medical kidnapping victim Alyssa Gilderhus at the Mayo Clinic. But was it actually  
→ InfoWars and FreeMartyG which publicly shamed CNN into doing this real journalism? Cue the  
→ Mission Impossible theme music for this one...\n\nThis mission, as we accepted it, began  
→ more than a year ago during the baby Charlie Gard medical kidnapping scandal in the UK and  
→ we thought that it had ended with an apparently unsuccessf'
```

Doing a LDA in Python

We will use gensim ([Rehurek10softwareframework](#)) for this (make sure you have version >4.0)

Let us assume you have a list of lists of documents called `texts`:

1

```
print(texts[0][:115])
```

which looks something like:

1

```
'Stop the presses: CNN covered some actual news yesterday when it reported on the story of  
→ medical kidnapping victim Alyssa Gilderhus at the Mayo Clinic. But was it actually  
→ InfoWars and FreeMartyG which publicly shamed CNN into doing this real journalism? Cue the  
→ Mission Impossible theme music for this one...\n\nThis mission, as we accepted it, began  
→ more than a year ago during the baby Charlie Gard medical kidnapping scandal in the UK and  
→ we thought that it had ended with an apparently unsuccessf'
```

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,

Preprocessing

Your preprocessing steps and feature engineering decisions *largely* affect your topics

1. You can apply 'manual' preprocessing steps ...
2. ... In isolation or combination with for example tfidf transformations

```
1 texts_clean = [text.lower() for text in texts]
2 texts_clean=[" ".join(text.split()) for text in texts_clean] #remove dubble spaces
3 texts_clean = [" ".join([l for l in text if l not in punctuation]) for text in texts_clean]
  ↳ #remove punctuation
4 texts_clean[0][:500]
```

which looks something like:

```
1 'stop the presses cnn covered some actual news yesterday when it reported on the story of
  ↳ medical kidnapping victim alyssa gilderhus at the mayo clinic but was it actually infowars
  ↳ and freemartyg which publicly shamed cnn into doing this real journalism cue the mission
  ↳ impossible theme music for this one this mission as we accepted it began more than a year
  ↳ ago during the baby charlie gard medical kidnapping scandal in the uk and we thought that
```

Preprocessing

Your preprocessing steps and feature engineering decisions *largely* affect your topics

1. You can apply 'manual' preprocessing steps ...
2. ... In isolation or combination with for example tfidf transformations

```
1 texts_clean = [text.lower() for text in texts]
2 texts_clean=[" ".join(text.split()) for text in texts_clean] #remove dubble spaces
3 texts_clean = [" ".join([l for l in text if l not in punctuation]) for text in texts_clean]
  ↳ #remove punctuation
4 texts_clean[0][:500]
```

which looks something like:

```
1 'stop the presses cnn covered some actual news yesterday when it reported on the story of
  ↳ medical kidnapping victim alyssa gilderhus at the mayo clinic but was it actually infowars
  ↳ and freemartyg which publicly shamed cnn into doing this real journalism cue the mission
  ↳ impossible theme music for this one this mission as we accepted it began more than a year
  ↳ ago during the baby charlie gard medical kidnapping scandal in the uk and we thought that
```

Preprocessing

Your preprocessing steps and feature engineering decisions *largely* affect your topics

1. You can apply 'manual' preprocessing steps ...
2. ... In isolation or combination with for example tfidf transformations

```
1 texts_clean = [text.lower() for text in texts]
2 texts_clean=[" ".join(text.split()) for text in texts_clean] #remove dubble spaces
3 texts_clean = [" ".join([l for l in text if l not in punctuation]) for text in texts_clean]
  ↳ #remove punctuation
4 texts_clean[0][:500]
```

which looks something like:

```
1 'stop the presses cnn covered some actual news yesterday when it reported on the story of
  ↳ medical kidnapping victim alyssa gilderhus at the mayo clinic but was it actually infowars
  ↳ and freemartyg which publicly shamed cnn into doing this real journalism cue the mission
  ↳ impossible theme music for this one this mission as we accepted it began more than a year
  ↳ ago during the baby charlie gard medical kidnapping scandal in the uk and we thought that
```

Preprocessing

Your preprocessing steps and feature engineering decisions *largely* affect your topics

1. You can apply 'manual' preprocessing steps ...
2. ... In isolation or combination with for example tfidf transformations

```
1 texts_clean = [text.lower() for text in texts]
2 texts_clean=[" ".join(text.split()) for text in texts_clean] #remove dubble spaces
3 texts_clean = [" ".join([l for l in text if l not in punctuation]) for text in texts_clean]
  ↳ #remove punctuation
4 texts_clean[0][:500]
```

which looks something like:

```
1 'stop the presses cnn covered some actual news yesterday when it reported on the story of
↳ medical kidnapping victim alyssa gilderhus at the mayo clinic but was it actually infowars
↳ and freemartyg which publicly shamed cnn into doing this real journalism cue the mission
↳ impossible theme music for this one this mission as we accepted it began more than a year
↳ ago during the baby charlie gard medical kidnapping scandal in the uk and we thought that
```

Preprocessing

Without *stopword removal*, *tfidf* transformation and/or *pruning*, you topics will not be very informative.

```
1  mystopwords = set(stopwords.words('english')) # use default NLTK stopwords list;  
    ↳ alternatively:  
2  # mystopwords = set(open('mystopwordfile.txt').readlines()) #read stopwords list from a  
    ↳ textfile with one stopwords per line  
3  texts_clean = [" ".join(word for word in text.split() if word not in mystopwords) for text  
    ↳ in texts_clean]  
4  texts_clean[0][:500]
```

which looks something like:

```
1  'stop presses cnn covered actual news yesterday reported story medical kidnapping victim  
    ↳ alyssa gilderhus mayo clinic actually infowars freemartyg publicly shamed cnn real  
    ↳ journalism cue mission impossible theme music one mission accepted began year ago baby  
    ↳ charlie gard medical kidnapping scandal uk thought ended apparently unsuccessful april  
    ↳ fools joke cnn sure many recall charlie gard infant rare form otherwise notsorare  
    ↳ condition mitochondrial disease story went viral made international news '
```


Preprocessing

Without *stopword removal*, *tfidf* transformation and/or *pruning*, you topics will not be very informative.

```
1  mystopwords = set(stopwords.words('english')) # use default NLTK stopwords list;  
    ↳ alternatively:  
2  # mystopwords = set(open('mystopwordfile.txt').readlines()) #read stopwords list from a  
    ↳ textfile with one stopwords per line  
3  texts_clean = [" ".join(word for word in text.split() if word not in mystopwords) for text  
    ↳ in texts_clean]  
4  texts_clean[0][:500]
```

which looks something like:

```
1  'stop presses cnn covered actual news yesterday reported story medical kidnapping victim  
    ↳ alyssa gilderhus mayo clinic actually infowars freemartyg publicly shamed cnn real  
    ↳ journalism cue mission impossible theme music one mission accepted began year ago baby  
    ↳ charlie gard medical kidnapping scandal uk thought ended apparently unsuccessful april  
    ↳ fools joke cnn sure many recall charlie gard infant rare form otherwise notsorare  
    ↳ condition mitochondrial disease story went viral made international news '
```

Preprocessing

Without *stopword removal*, *tfidf* transformation and/or *pruning*, you topics will not be very informative.

```
1  mystopwords = set(stopwords.words('english')) # use default NLTK stopwords list;  
    ↳ alternatively:  
2  # mystopwords = set(open('mystopwordfile.txt').readlines()) #read stopwords list from a  
    ↳ textfile with one stopwords per line  
3  texts_clean = [" ".join(word for word in text.split() if word not in mystopwords) for text  
    ↳ in texts_clean]  
4  texts_clean[0][:500]
```

which looks something like:

```
1  'stop presses cnn covered actual news yesterday reported story medical kidnapping victim  
    ↳ alyssa gilderhus mayo clinic actually infowars freemartyg publicly shamed cnn real  
    ↳ journalism cue mission impossible theme music one mission accepted began year ago baby  
    ↳ charlie gard medical kidnapping scandal uk thought ended apparently unsuccessful april  
    ↳ fools joke cnn sure many recall charlie gard infant rare form otherwise notsorare  
    ↳ condition mitochondrial disease story went viral made international news '
```

Tokenization

gensim expects a list of words (hence: tokenize your corpus)

```
1  tokenized_texts_clean = [TreebankWordTokenizer().tokenize(text) for text in texts_clean ] #  
    ↪ tokenize texts; convert all strings to a list of tokens  
2  tokenized_texts_clean[0][:500]
```

which looks something like:

```
1  ['stop',  
2  'presses',  
3  'cnn',  
4  'covered',  
5  'actual',  
6  'news',  
7  'yesterday',  
8  'reported',  
9  'story',  
10 ..
```

Tokenization

gensim expects a list of words (hence: tokenize your corpus)

```
1  tokenized_texts_clean = [TreebankWordTokenizer().tokenize(text) for text in texts_clean ] #  
    ↪ tokenize texts; convert all strings to a list of tokens  
2  tokenized_texts_clean[0][:500]
```

which looks something like:

```
1  ['stop',  
2  'presses',  
3  'cnn',  
4  'covered',  
5  'actual',  
6  'news',  
7  'yesterday',  
8  'reported',  
9  'story',  
10 ..
```

Tokenization

gensim expects a list of words (hence: tokenize your corpus)

```
1  tokenized_texts_clean = [TreebankWordTokenizer().tokenize(text) for text in texts_clean ] #  
    ↪ tokenize texts; convert all strings to a list of tokens  
2  tokenized_texts_clean[0][:500]
```

which looks something like:

```
1  ['stop',  
2  'presses',  
3  'cnn',  
4  'covered',  
5  'actual',  
6  'news',  
7  'yesterday',  
8  'reported',  
9  'story',  
10 ..
```

LDA implementation

LDA implementation

```

1 raw_m1 = tokenized_texts_clean
2
3 # assign a token_id to each word
4 id2word_m1 = corpora.Dictionary(raw_m1)
5 # represent each text by (token_id, token_count) tuples
6 ldacorporus_m1 = [id2word_m1.doc2bow(text) for text in raw_m1]
7
8 #estimate the model
9 lda_m1 = models.LdaModel(ldacorporus_m1, id2word=id2word_m1, num_topics=10)
10 lda_m1.print_topics()

```

```

1 [(0, '0.015*"trump" + 0.012*"said" + 0.006*"president" + 0.006*"people" + 0.004*"cnn" +
↳ 0.004*"us" + 0.004*"house" + 0.004*"news" + 0.003*"also" + 0.003*"twitter"'),
2 (1, '0.010*"said" + 0.008*"trump" + 0.004*"one" + 0.004*"people" + 0.004*"us" +
↳ 0.004*"president" + 0.004*"would" + 0.003*"media" + 0.003*"also" + 0.003*"new"'),
3 (2, '0.011*"trump" + 0.009*"said" + 0.007*"president" + 0.005*"would" + 0.004*"people" +
↳ 0.004*"us" + 0.003*"also" + 0.003*"like" + 0.003*"news" + 0.003*"state"'),
4 (3, '0.010*"trump" + 0.006*"president" + 0.005*"said" + 0.004*"would" + 0.004*"us" +
↳ 0.003*"also" + 0.003*"people" + 0.003*"media" + 0.003*"news" + 0.003*"one"'),

```

LDA implementation

LDA implementation

```

1 raw_m1 = tokenized_texts_clean
2
3 # assign a token_id to each word
4 id2word_m1 = corpora.Dictionary(raw_m1)
5 # represent each text by (token_id, token_count) tuples
6 ldacorporus_m1 = [id2word_m1.doc2bow(text) for text in raw_m1]
7
8 #estimate the model
9 lda_m1 = models.LdaModel(ldacorporus_m1, id2word=id2word_m1, num_topics=10)
10 lda_m1.print_topics()
```

```

1 [(0, '0.015*"trump" + 0.012*"said" + 0.006*"president" + 0.006*"people" + 0.004*"cnn" +
↳ 0.004*"us" + 0.004*"house" + 0.004*"news" + 0.003*"also" + 0.003*"twitter"'),
2 (1, '0.010*"said" + 0.008*"trump" + 0.004*"one" + 0.004*"people" + 0.004*"us" +
↳ 0.004*"president" + 0.004*"would" + 0.003*"media" + 0.003*"also" + 0.003*"new"'),
3 (2, '0.011*"trump" + 0.009*"said" + 0.007*"president" + 0.005*"would" + 0.004*"people" +
↳ 0.004*"us" + 0.003*"also" + 0.003*"like" + 0.003*"news" + 0.003*"state"'),
4 (3, '0.010*"trump" + 0.006*"president" + 0.005*"said" + 0.004*"would" + 0.004*"us" +
↳ 0.003*"also" + 0.003*"people" + 0.003*"media" + 0.003*"news" + 0.003*"one"'),
```

LDA implementation

LDA implementation

```

1 raw_m1 = tokenized_texts_clean
2
3 # assign a token_id to each word
4 id2word_m1 = corpora.Dictionary(raw_m1)
5 # represent each text by (token_id, token_count) tuples
6 ldacorporus_m1 = [id2word_m1.doc2bow(text) for text in raw_m1]
7
8 #estimate the model
9 lda_m1 = models.LdaModel(ldacorporus_m1, id2word=id2word_m1, num_topics=10)
10 lda_m1.print_topics()

```

```

1 [(0, '0.015*"trump" + 0.012*"said" + 0.006*"president" + 0.006*"people" + 0.004*"cnn" +
↳ 0.004*"us" + 0.004*"house" + 0.004*"news" + 0.003*"also" + 0.003*"twitter"'),
2 (1, '0.010*"said" + 0.008*"trump" + 0.004*"one" + 0.004*"people" + 0.004*"us" +
↳ 0.004*"president" + 0.004*"would" + 0.003*"media" + 0.003*"also" + 0.003*"new"'),
3 (2, '0.011*"trump" + 0.009*"said" + 0.007*"president" + 0.005*"would" + 0.004*"people" +
↳ 0.004*"us" + 0.003*"also" + 0.003*"like" + 0.003*"news" + 0.003*"state"'),
4 (3, '0.010*"trump" + 0.006*"president" + 0.005*"said" + 0.004*"would" + 0.004*"us" +
↳ 0.003*"also" + 0.003*"people" + 0.003*"media" + 0.003*"news" + 0.003*"one"'),

```


Visualization with pyldavis

```
1 import pyLDavis
2 import pyLDavis.gensim_models as gensimvis
3 # first estimate gensim model, then:
4 vis_data = gensimvis.prepare(lda_m1,ldacorporus_m1,id2word_m1)
5 pyLDavis.display(vis_data)
```

Unsupervised Machine Learning for Text Classification

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 of 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

Evaluation metrics

Qualitative: human judgement

Observation and interpretation based: observe the top N words in your topic, and evaluate the quality of the coherence of the topic. Can you identify words that do not belong to a topic?

Quantitative: 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()`)

Evaluation metrics

Qualitative: human judgement

Observation and interpretation based: observe the top N words in your topic, and evaluate the quality of the coherence of the topic. Can you identify words that do not belong to a topic?

Quantitative: 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()`)

So, how do we do this?

- Estimate multiple models, store the metrics for each model, and then compare them (numerically, or by plotting)
- Idea: We select some candidate models, and then look whether they can be interpreted.
- But what can we tune?

So, how do we do this?

- Estimate multiple models, store the metrics for each model, and then compare them (numerically, or by plotting)
- Idea: We select some candidate models, and then look whether they can be interpreted.
- But what can we tune?

So, how do we do this?

- Estimate multiple models, store the metrics for each model, and then compare them (numerically, or by plotting)
- Idea: We select some candidate models, and then look whether they can be interpreted.
- But what can we tune?

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 peculiarities instead of finding general patterns
- There is no inherent ordering of topics
- 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 α : how sparse should the document-topic distribution θ be?

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

1

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

Choosing α : how sparse should the document-topic distribution θ be?

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

1

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

Unsupervised Machine Learning for Text Classification

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).



Any questions?

Next steps

Take-home exam: you have time until
Thursday 11th, end of the day

An example notebook with code for running
LDA models can be found here:

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

References



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. <https://doi.org/10.1080/21670811.2015.1096598>