# Tutorial: Topic Modeling

Niels Scholten | JM2050-M-6

TU/e EINDHOVEN UNIVERSITY OF TECHNOLOGY

TILBURG UNIVERSITY

JADS Jheronimus Academy of Data Science

# Today

Learning Goals:

- Familiarize yourself with the training processs of different topic models
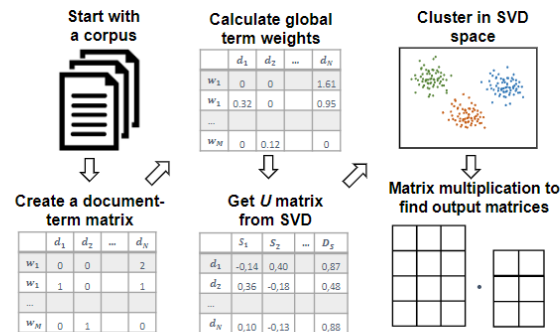
- Familiarize yourself with topic evaluations

# FuzzyTM

[FuzzyTM](#)
A python library for training fuzzy topic models and creating topic embeddings for downstream tasks.

# FuzzyTM

- Algorithms
  - FLSA
  - **FLSA-W (For the assignment)**
  - FLSA-E
  - FLSA-V

The algorithms differ in the space in which they cluster

- Global-term weighting
  - Entropy
  - IDF
  - **Normal (Default)**
  - PorbIDF

- Clustering Algorithm
  - **FCM (Default)**
  - Gustafson-Kessel
  - FST-PSO

# Other useful methods

- Evaluation scores:
  - get_coherence_score()
  - get_diversity_score()
  - get_interpretability_score()

- Descriptive information:
  - get_vocabulary()
  - get_vocabulary_size()

- Embeddings:
  - get_topic_embedding()

# Calculate metrics for external topics!

```
flsaW.get_coherence_score(
                data,
                [['fax','modem','card',],
                [fax','disk','write',],],
                )


flsaW.get_diversity_score(
                data,
                [['fax','modem','card',],
                [fax','disk','write',],],
                )
```

# Example

# Load data

```python
from nltk.tokenize import word_tokenize
import re
from unidecode import unidecode
import pandas as pd
from nltk.corpus import stopwords
stop_words = set(stopwords.words('english'))
from collections import Counter

FILE_PATH = <your_path>
data = pd.read_csv(
            FILE_PATH,
            usecols=['Text'],
            )['Text'].tolist()
```

# Preprocess data

```python
def preprocess_texts(texts, n=0):
    processed_texts = []
        for text in texts:
            # lowercasing, keep text only, remove accents, tokenization
            tokens = [word for word in word_tokenize(re.sub(r'[^a-zA-Z\s]', '', unidecode(text.lower())))]
            # stopword removal
            tokens = [token for token in tokens if token not in stopwords.words('english')]
            processed_texts.append(tokens)

    # remove top-n% and bottom-n% words
    if n > 0:
        word_freq = Counter([word for sentence in processed_texts for word in sentence])
        top_n = set([word for word, _ in word_freq.most_common(int(n/100*len(word_freq)))])
        bottom_n = set([word for word, _ in word_freq.most_common()[:-int(n/100*len(word_freq))-1:-1]])
        processed_texts = [[word for word in sentence if word not in top_n and word not in bottom_n] for sentence in
processed_texts]

    return processed_texts
```

# Filter data and tokenize

**Filter:**

```python
def filter_word_from_corpus(data, words):
    # Ensure words is a list, even if a single string is passed
    if isinstance(words, str):
        words = [words]
    # Filter words from data
        filtered_data = [[token for token in row if token not in words] for row in data]
    return filtered_data
```

**Tokenize:**

```python
s = "This is a sample string."
tokens = s.split() # Splits on whitespace by default
print(tokens) # ['This', 'is', 'a', 'sample', 'string.']
```

# Training FLSA-W with FuzzyTM

```python
from FuzzyTM import FLSA_W
flsaW = FLSA_W(
    input_file = data,
    num_topics=10,
    num_words=10,
    )

flsaW.get_vocabulary_size()

pwgt, ptgd = flsaW.get_matrices() # THIS TRAINS THE MODEL

print(flsaW.show_topics())
for topic in flsaW.show_topics(representation='words'):
    print(topic)

print(flsaW.get_coherence_score())
print(flsaW.get_diversity_score())
print(flsaW.get_interpretability_score())
```

# Training an LDA Model with Gensim

```python
import gensim
from gensim import corpora
from gensim.models.ldamodel import LdaModel
from gensim.models import CoherenceModel
dictionary = corpora.Dictionary(data)

# Convert the list of documents (corpus) into Document Term Matrix using the dictionary prepared above
doc_term_matrix = [dictionary.doc2bow(doc) for doc in data]

# Train the LDA model
lda_model = LdaModel(
    doc_term_matrix,
    num_topics=3,
    id2word=dictionary,
    )

topics = lda_model.print_topics()
# Print the topics
for topic in lda_model.print_topics():
    print(topic)

#Get coherence score
print(CoherenceModel(model=lda_model, texts=data, dictionary=dictionary, coherence='c_v').get_coherence())
```

# Training BERTopic

- If you do not use a GPU, this could take a while. Google Colab could speed it up

```python
from bertopic import BERTopic
docs = [' '.join(doc) for doc in data]

# Create the model (uses DistilBERT by default)
bert_topic = BERTopic()

# Train the model and transform your data into topics
topic_assigned_to_doc, _ = bert_topic.fit_transform(docs)

topic_matrix = bert_topic.get_topic_info()
```
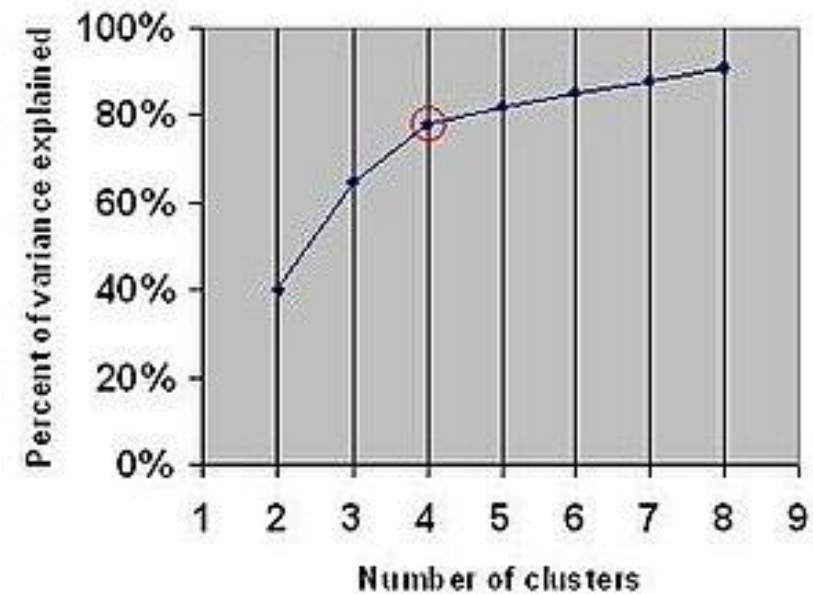
# Elbow Plot

- Heuristic often used in clustering
- Can be used in FLSA-W to determine the number of clusters
- Find the number where the curve "breaks"

- However, you might not a nice curve like in the image

- For FLSA-W, use topic coherence instead for *Percent of variance explained.*
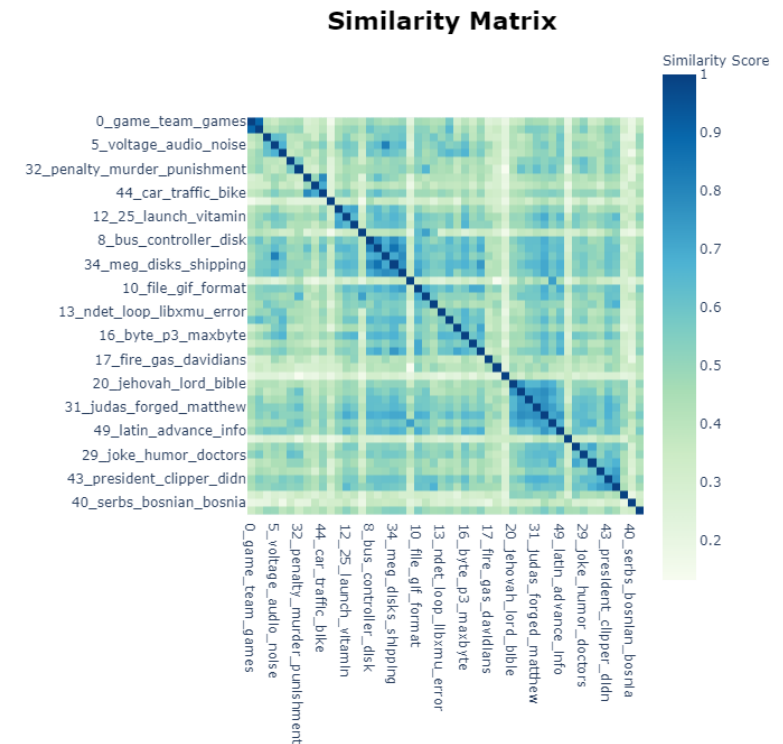
# Topic Similarity Matrix

- Nice feature in Bertopic
- Calculates the cosine similarity between topic embeddings

`Topic_model.visualize_heatmap()`



Similarity Matrix

# Exercise

Niels Scholten | JM2050-M-6

JADS Jheronimus Academy of Data Science

# For FLSA-W and LDA

Train topic models with both algorithms.
- Use an iterative process in which you:
  - train a topic model,
  - remove useless words from the corpus ,
  - retrain your model.

- How many iterations did it take to find your topics?
- What is the effect on the coherence score?

A few considerations:
- Setting the number of words per topics: Which number seems most intuitive?
- Setting the number of topics: How did you select this number?

# BERTopic

For the first ten documents in `topic_assigned_to_doc`, compare the produced topic to the document. Do you agree this document reflects the topic well?

What do you think -1 means?

How many words per topic do you have? What is the effect on the topic quality if you work with more topics?

# All models

- Train all models with the same number of words per topic

- Calculate the coherence scores for all three models
  - How do the models compare?
  - Does your personal judgement align with the scores?

# Interpreting Topics

- Ask ChatGPT to give names to the topics

- Are these descriptions good?

- How do you prompt ChatGPT to give good topic names and descriptions?