



# Introduction to Topic Modeling

## Saeed Roshani (Ph.D.)





The goal of this study is to develop a technology analysis of the evolution of some critical quantum technologies to improve technological forecasting. Methodology applies S-shaped curves based on patent data to analyze the evolutionary phases of quantum technologies over the course of time.

"This is a sample"

↓  
Tokenization

↓  
"This" "is" "a" "sample"

The goal of this study is to develop a technology analysis of the evolution of some critical quantum technologies to improve technological forecasting. Methodology applies S-shaped curves based on patent data to analyze the evolutionary phases of quantum technologies over the course of time.

```
# Download necessary data and models  
import nltk  
nltk.download('stopwords')  
nlp = spacy.load('en_core_web_sm')
```

The goal of this study is to develop a technology analysis of the evolution of some critical quantum technologies to improve technological forecasting. Methodology applies S-shaped curves based on patent data to analyze the evolutionary phases of quantum technologies over the course of time.

```
# Define stop words
stop_words = stopwords.words('english')
stop_words.extend(['research', 'use', 'new', 'article', 'study', 'model', 'define', 'background', 'examine', 'paper',
                  'examination', 'investigate', 'aim', 'goal', 'objective', 'US', 'china', 'result',
                  'measure', 'measurement', 'quality', 'analysis', 'approach', 'method', 'propose', 'investigate'])
```

evolution  
technological forecasting.  
patent  
technologies

technology analysis  
quantum technologies  
S-shaped curves  
evolutionary phases  
quantum

.

```
def make_bigrams(texts):
    bigram = Phrases(texts, min_count=10, threshold=50)
    bigram_phraser = Phraser(bigram)
    trigram = Phrases(bigram_phraser[texts], min_count=10, threshold=50)
    trigram_phraser = Phraser(trigram)
    return [trigram_phraser[bigram_phraser[doc]] for doc in texts]
```

evolution  
technological forecasting.  
patent  
technolog(y)ies

technology analysis  
quantum technolog(y)ies  
applies S-shaped curves  
evolutionary phases quantum  
.

```
def lemmatize(text, allowed_postags=['NOUN', 'VERB', 'ADJ', 'ADV']):  
    doc = nlp(text)  
    return [token.lemma_ for token in doc if token.pos_ in allowed_postags and token.lemma_ not in stop_words]
```

evolution  
technolog  
patent  
technolog

forecast .

quantum technolog  
evolutionary phase .

technolog analysi  
S-shape curve  
quantum

```
def stem(text):
    p = PorterStemmer()
    return [p.stem(token) for token in text]
```



```
import gensim

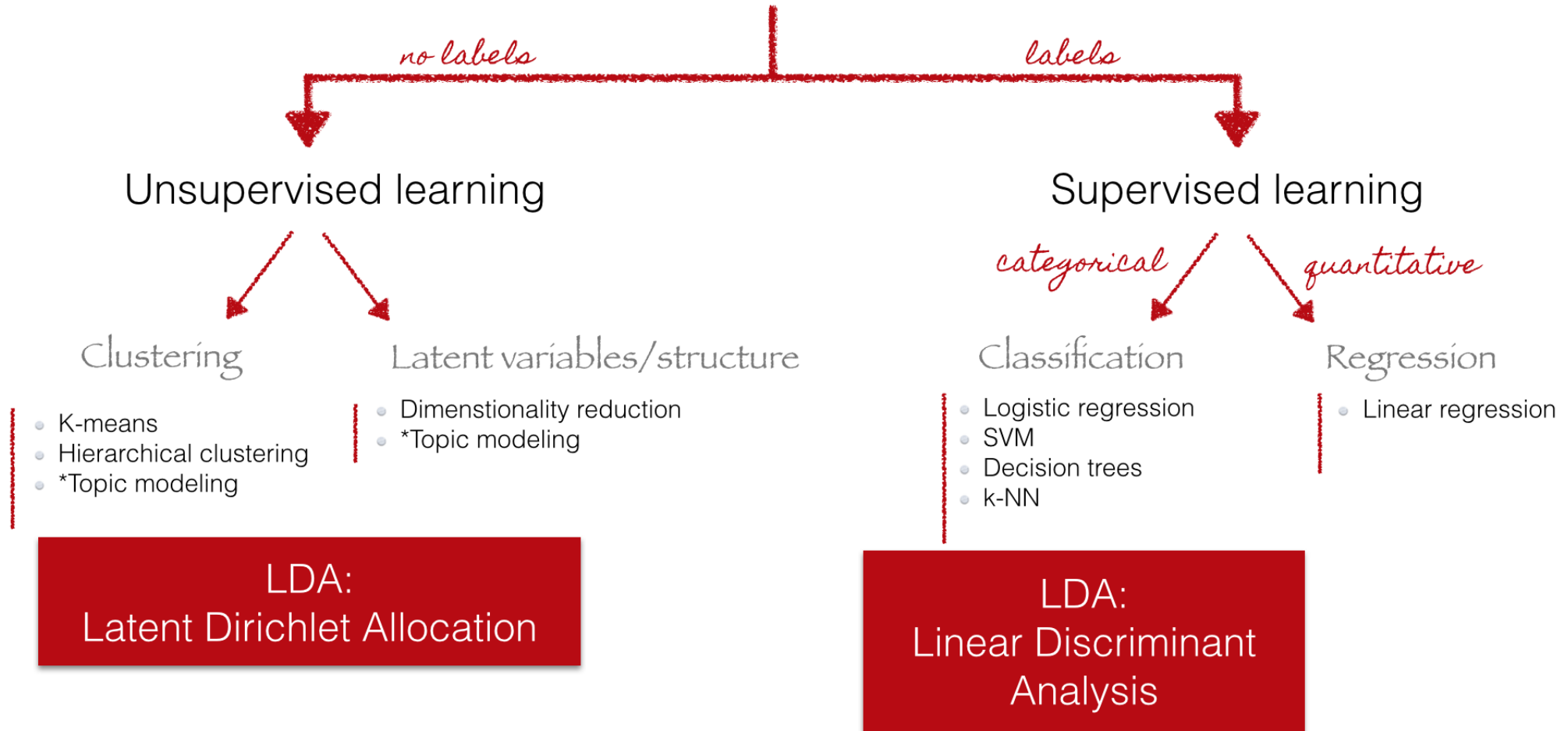
# Create a list of lists of the cleaned tokens
text = df['stemmed'].tolist()

# Create a dictionary of the cleaned tokens
dictionary = gensim.corpora.Dictionary(text)
print('Total Vocabulary Size:', len(dictionary))

# Filter out words that occur less than 5 documents, or more than 90% of the documents.
dictionary.filter_extremes(no_below=2, no_above=0.7)
print('Total Vocabulary Size after removing:', len(dictionary))

# Create a bag-of-words corpus of the cleaned tokens
corpus = [dictionary.doc2bow(t) for t in text]
```

# Machine Learning



# Definitions

A topic model is a type of statistical model for **discovering** the abstract "**topics**" that occur in a collection of **documents**.

Topic models are a suite of algorithms that uncover the **hidden thematic** structure in document collections. These algorithms help us develop new ways to search, browse and summarize large archives of texts.

Topic models provide a simple way to analyze large volumes of **unlabeled text**. A "topic" consists of a **cluster** of words that **frequently** occur together



# Latent Dirichlet Allocation (LDA)

Latent Dirichlet Allocation (LDA) is a **generative statistical model** used for topic modeling, which aims to identify the underlying themes or topics in a collection of documents. In LDA, a document is seen as a mixture of various **topics** and each topic is represented by a distribution over words in the vocabulary.

The goal of LDA is to infer the topic distribution in each document and the word distribution in each topic based on the observed words in the corpus. This is done through **Bayesian inference**, where the posterior distribution over the latent variables (i.e., the topic distribution in each document and the word distribution in each topic) is approximated using **Markov chain Monte Carlo** (MCMC).

By simulating this Markov chain for a large number of iterations, we can **approximate** the posterior distribution over the latent variables. This allows us to estimate the most likely topic and word distributions that generated the observed data, and to make inferences about the relationships between topics and words in the corpus.

---

# Latent Dirichlet Allocation

---

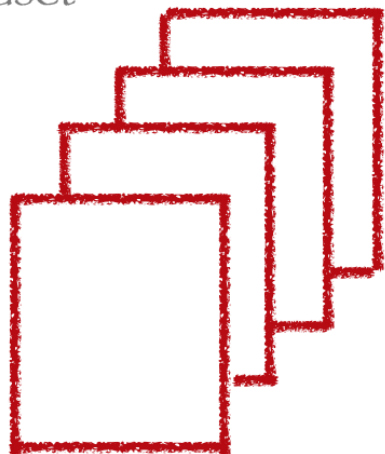
**David M. Blei, Andrew Y. Ng and Michael I. Jordan**  
University of California, Berkeley  
Berkeley, CA 94720

## Abstract

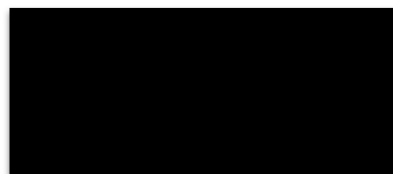


We propose a generative model for text and other collections of discrete data that generalizes or improves on several previous models including naive Bayes/unigram, mixture of unigrams [6], and Hofmann's aspect model, also known as probabilistic latent semantic indexing (pLSI) [3]. In the context of text modeling, our model posits that each document is generated as a mixture of topics, where the continuous-valued mixture proportions are distributed as a latent Dirichlet random variable. Inference and learning are carried out efficiently via variational algorithms. We present empirical results on applications of this model to problems in text modeling, collaborative filtering, and text classification.

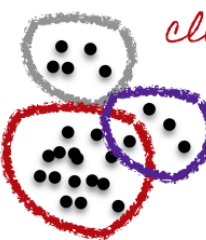
Dataset



*collection of text documents*



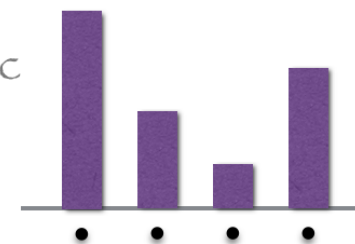
Topics



*cluster of words*

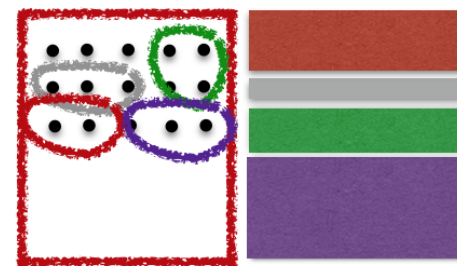
• words/tokens

Topic



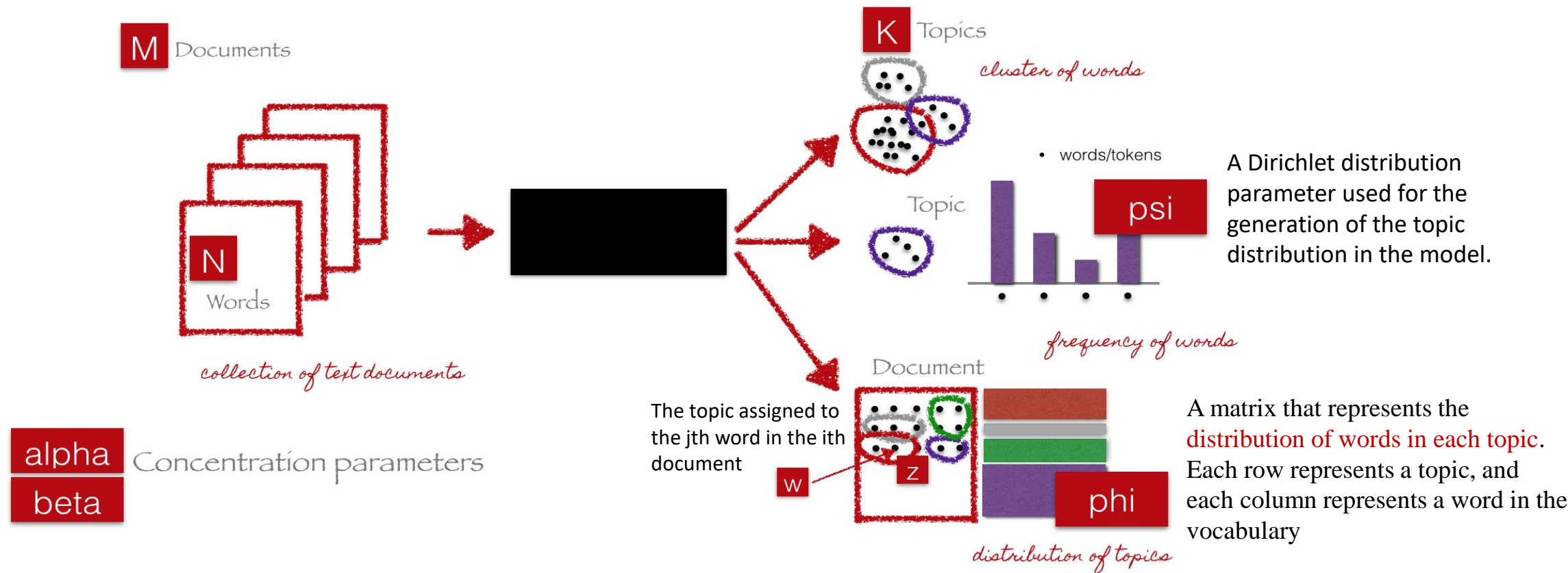
*frequency of words*

Document



*distribution of topics*





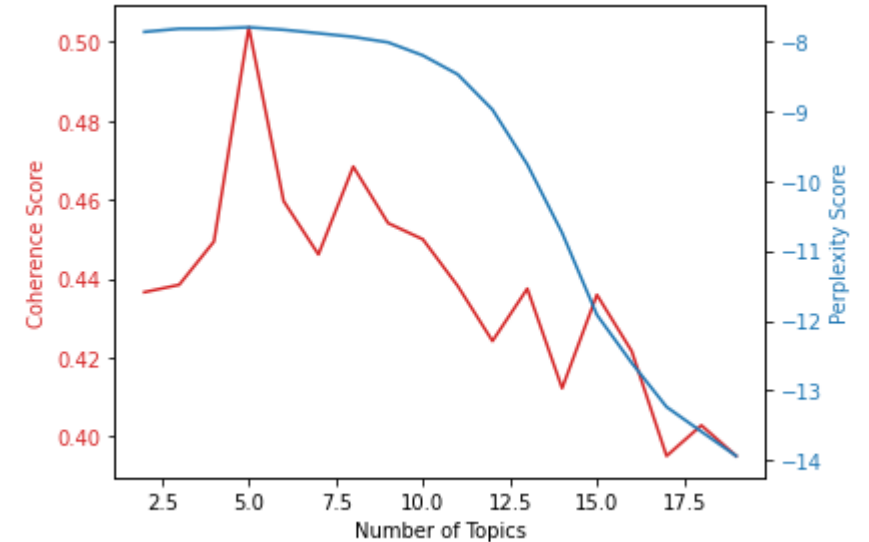
alpha: A hyperparameter that controls the distribution of **topics** in a document. A higher alpha value results in a document containing more topics.  
 Beta: A hyperparameter that controls the distribution of **words** in a topic. A higher beta value results in a topic containing more words.

# How to find optimal number of topics?

**Coherence score** and **perplexity** are two common evaluation metrics used in topic modeling.

Coherence score measures how well the words in a topic relate to each other and how distinct they are from the words in other topics.

Perplexity measures the ability of the model to predict new documents. It calculates how well the model can predict the probability of unseen words in a held-out test set. A lower perplexity score indicates that the model is better at predicting unseen data.



Both coherence score and perplexity are useful in evaluating the performance of a topic model. However, it's important to keep in mind that they don't necessarily reflect the real-world usefulness or relevance of the topics generated by the model. Therefore, it's always recommended to supplement quantitative evaluations with qualitative assessments by domain experts.

```
# Train an LDA model with the optimal number of topics
```

```
lda_model = LdaModel(corpus=corpus,
```

```
    id2word=dictionary, id2word: mapping of the vocabulary IDs to the actual words in the corpus
```

```
    num_topics=15,
```

```
    chunksize=100, specifies the number of documents to be used in each training chunk
```

```
    passes=10,
```

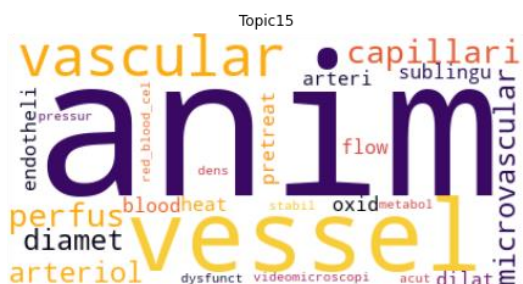
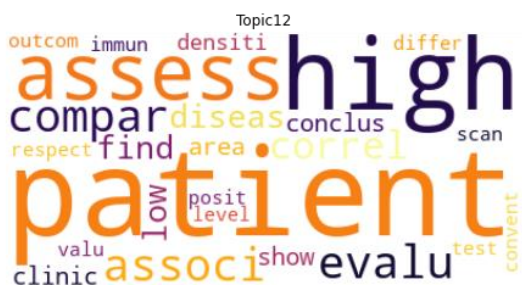
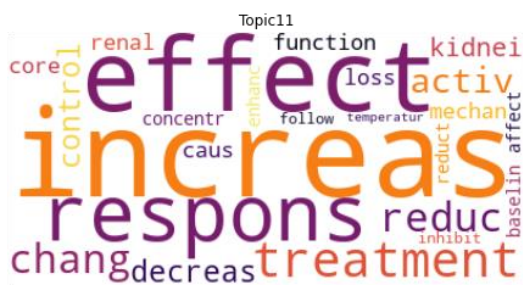
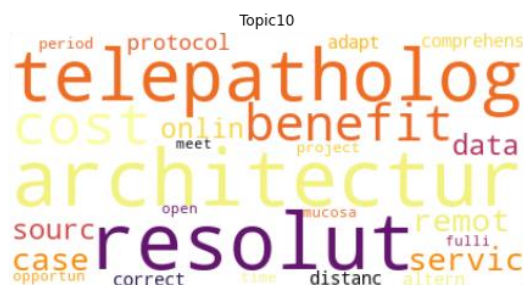
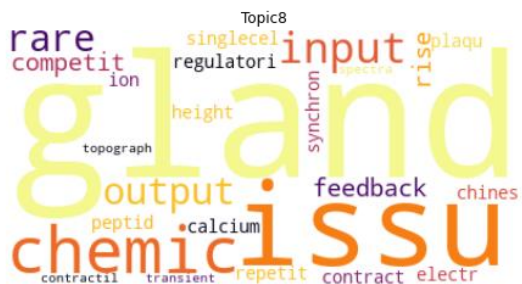
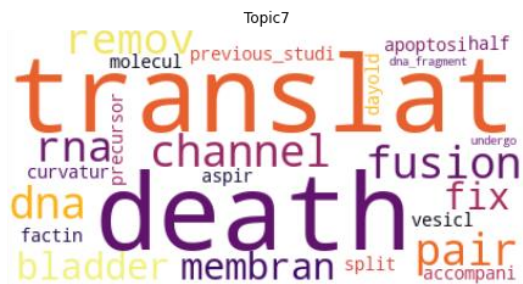
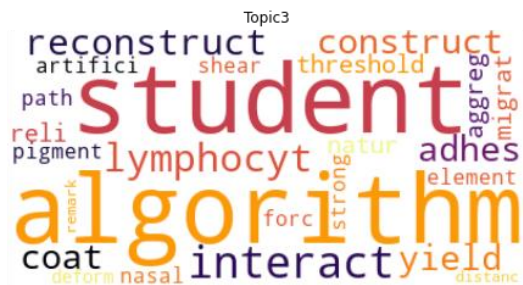
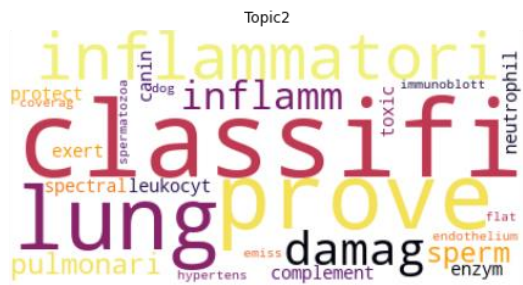
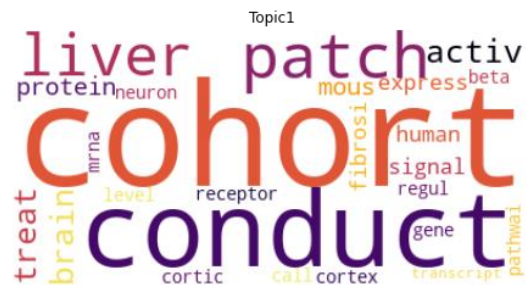
```
    alpha='auto', specifies the number of times the corpus will be processed during training
```

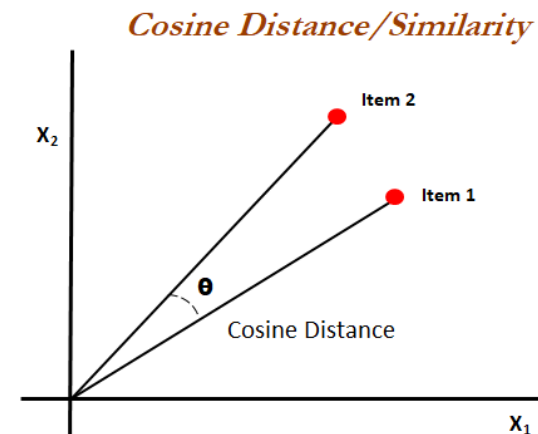
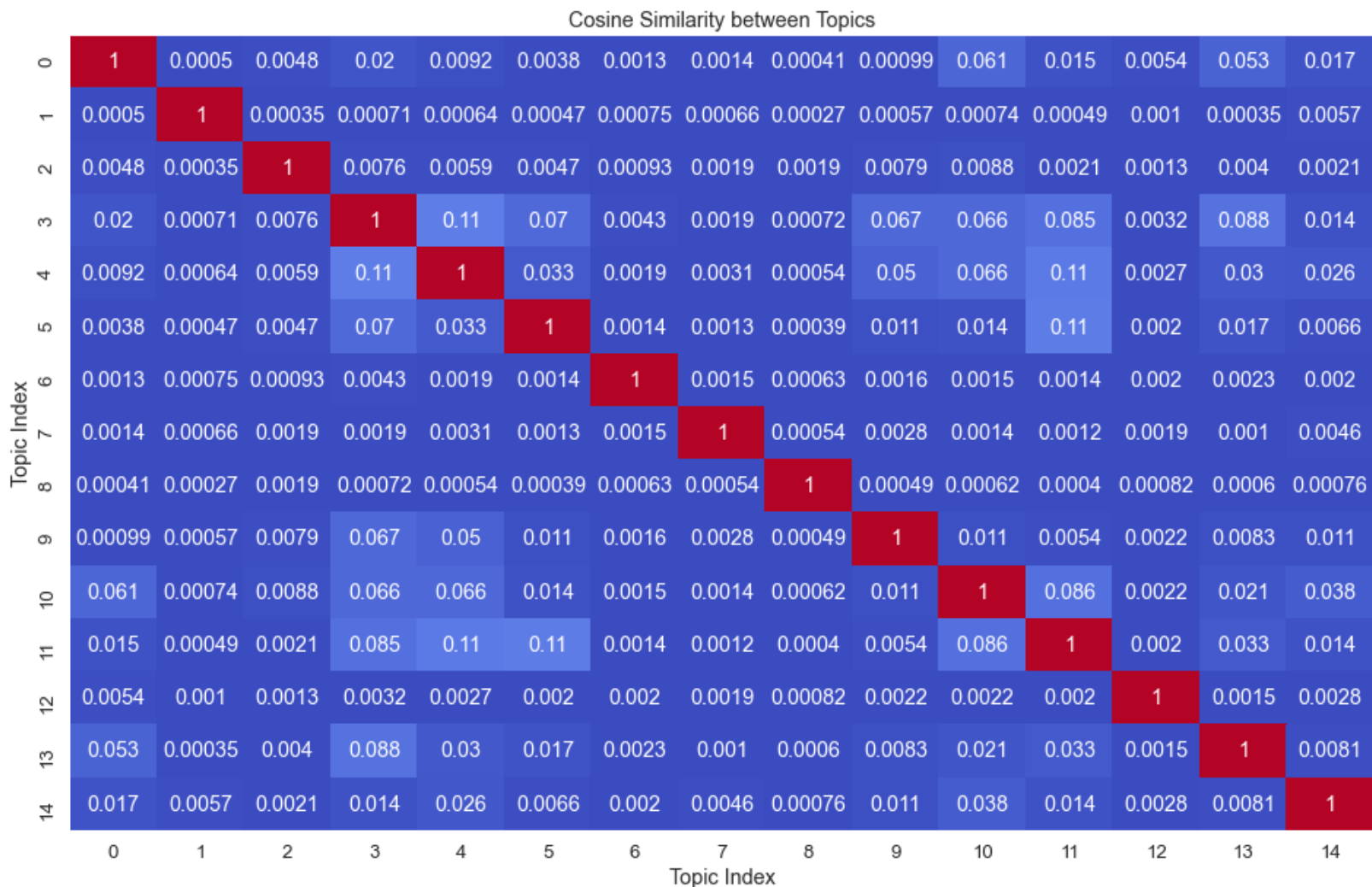
```
    per_word_topics=True)
```

```
for topic in lda_model.print_topics():
```

```
    print(topic)
```



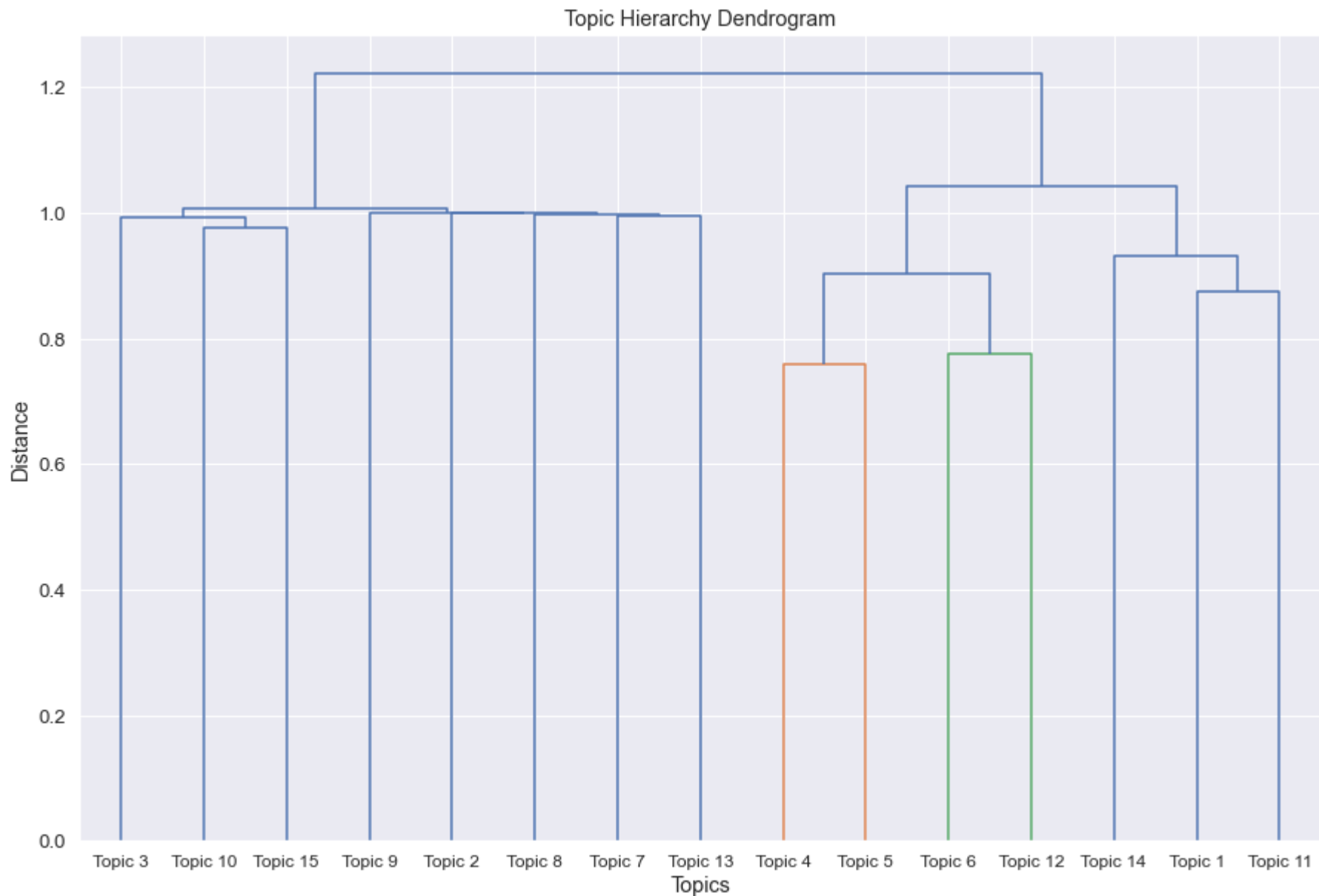




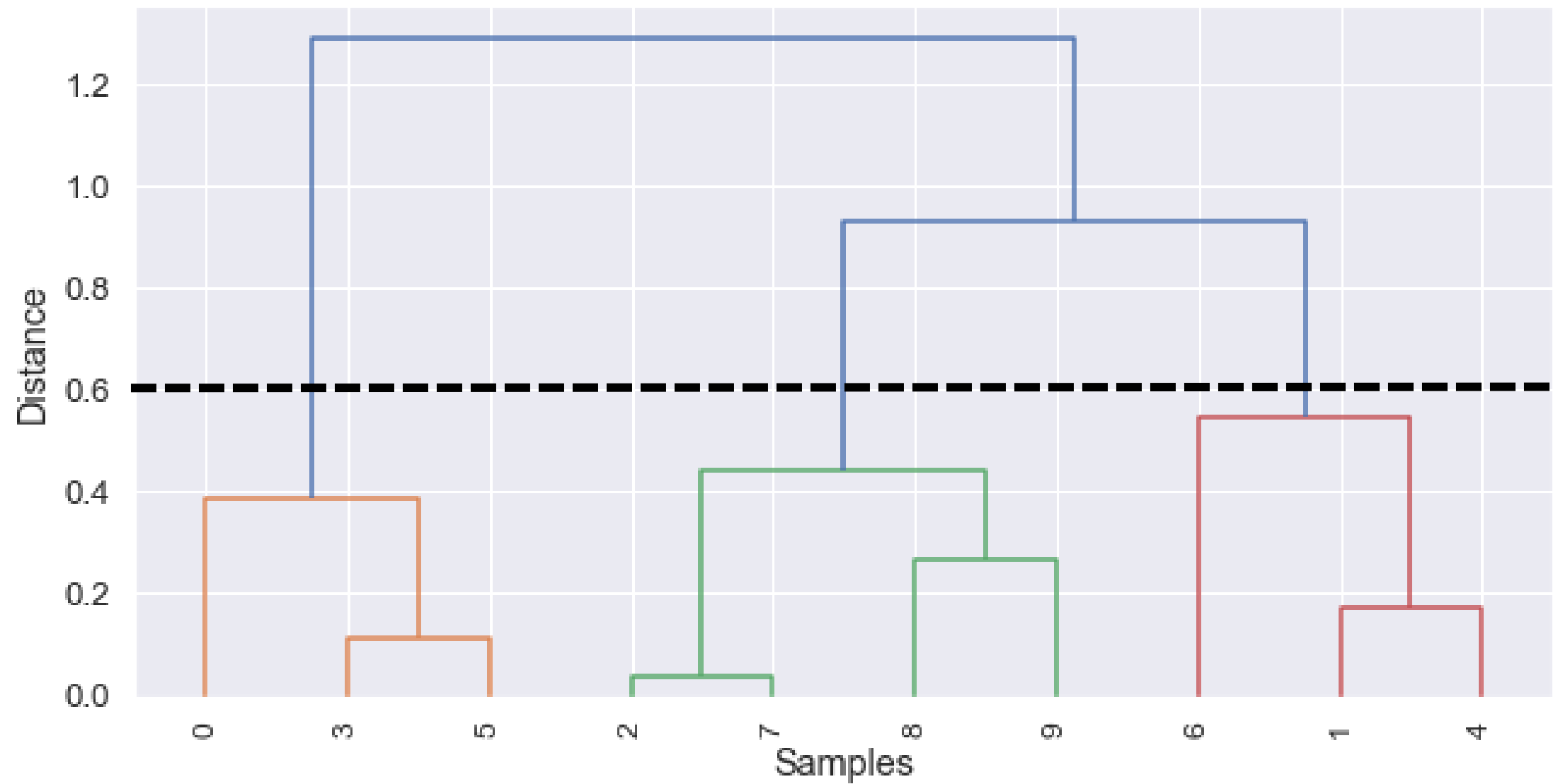
$$\cos(\theta) = \frac{\mathbf{A} \cdot \mathbf{B}}{\|\mathbf{A}\| \|\mathbf{B}\|} = \frac{\sum_{i=1}^n A_i B_i}{\sqrt{\sum_{i=1}^n A_i^2} \sqrt{\sum_{i=1}^n B_i^2}}$$

Topic 1: 0.3 \* word1 + 0.4 \* word2 + 0.3 \* word3

the value 0.3 for word 1 in topic 1 represents the probability of word 1 being generated from topic 1. In other words, if we randomly select a word from a document that belongs to topic 1, there is a 30% chance that the word will be word 1.

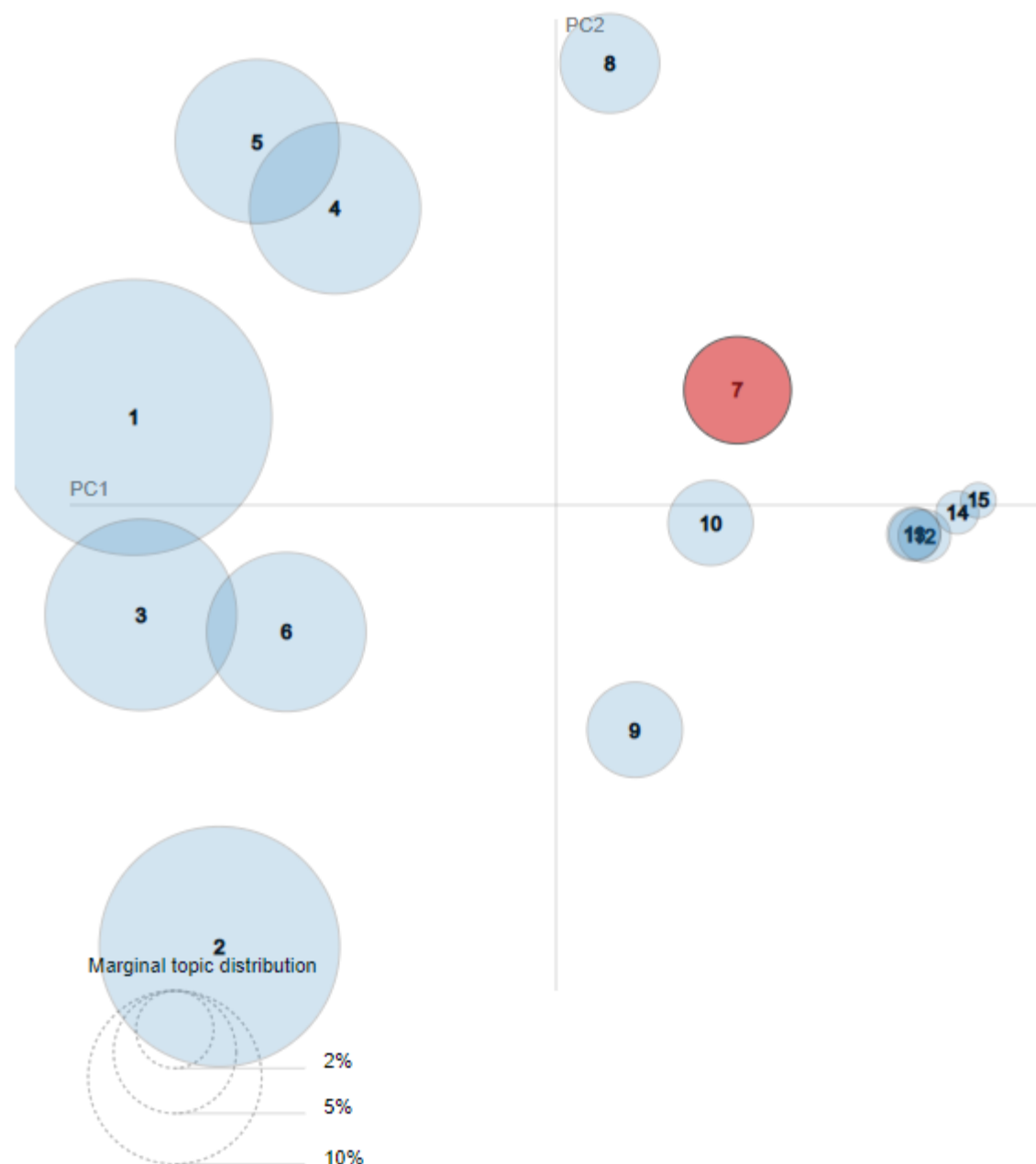


Dendrogram

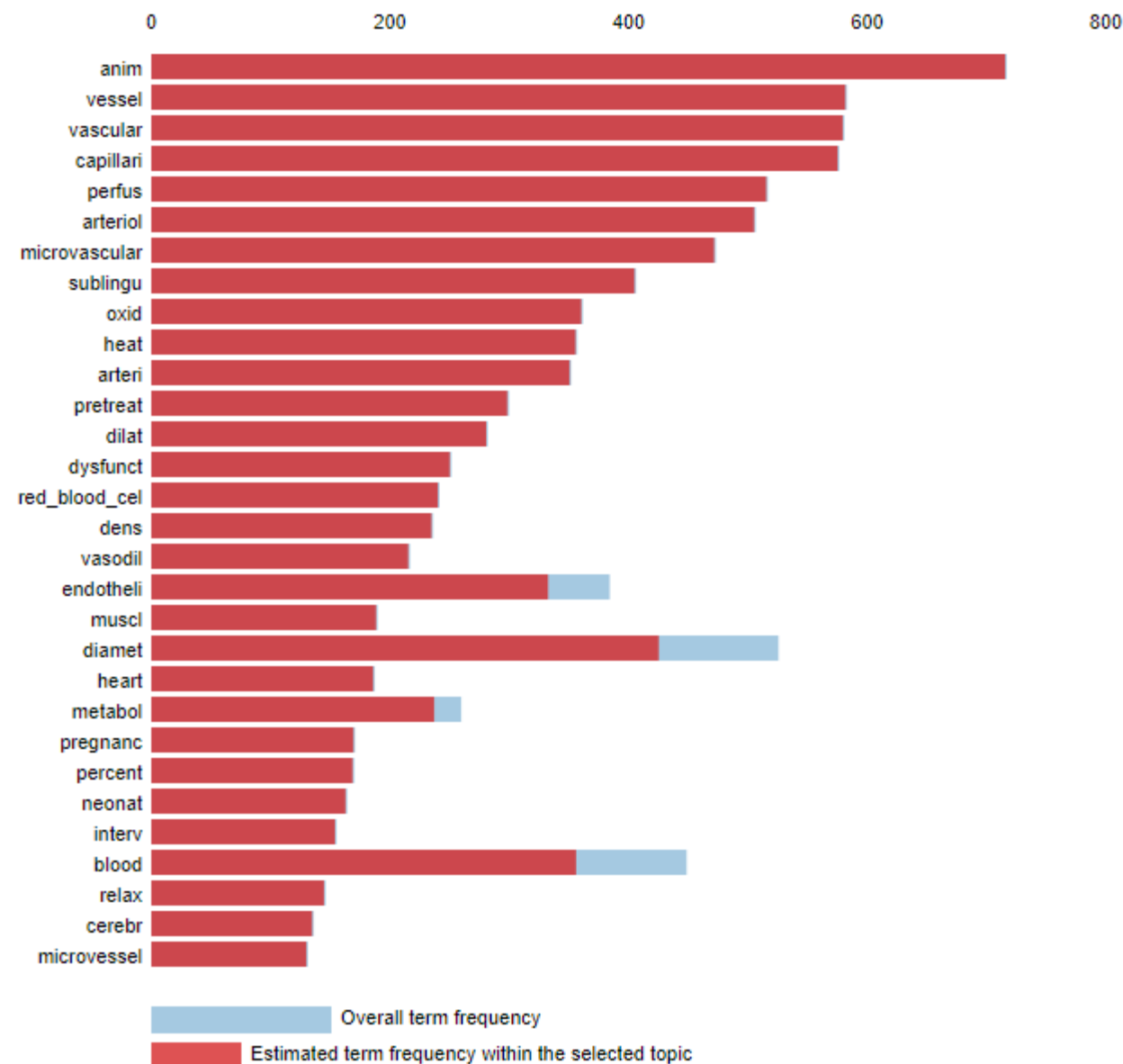




Intertopic Distance Map (via multidimensional scaling)



Top-30 Most Relevant Terms for Topic 7 (3.8% of tokens)



1.  $\text{saliency}(\text{term } w) = \text{frequency}(w) * [\sum_t p(t | w) * \log(p(t | w) / p(t))]$  for topics  $t$ ; see Chuang et. al (2012)

2.  $\text{relevance}(\text{term } w | \text{topic } t) = \lambda * p(w | t) + (1 - \lambda) * p(w | t) / p(w)$ ; see Sievert & Shirley (2014)