

Text Analysis: Health Misinformation

Sharon Jepkosgei

Christopher Womble

Kendall Goodland



Background and Literature

- 80% of people online use the internet to search for health information yet it's plagued with fake news.
- Many studies on fake news in other domains but few on health sector.
- Those in health domain are mainly COVID-19 specific.
- Our study – Understand health misinformation in a more general sense.



The slide features a dark blue background with decorative elements on the left. These include a large solid orange hexagon at the top, a smaller solid light blue hexagon at the bottom, and a hexagonal inset showing a collage of various data charts and graphs. A thin orange hexagonal outline is also visible, partially overlapping the light blue hexagon.

Research Questions

1. Does health misinformation exhibits distinct language patterns, with certain words appearing more than other?
2. Does health misinformation have high prevalence of certain topics?



Research Hypothesis

1. Health misinformation exhibits distinct linguistic patterns, with certain words recurring more frequently than others.
2. Health misinformation has distinct thematic patterns, with certain topics appearing more than others.

Data



Source

- HealthNewsReview.org.
- Health publication reviewed by experts



Sample Size

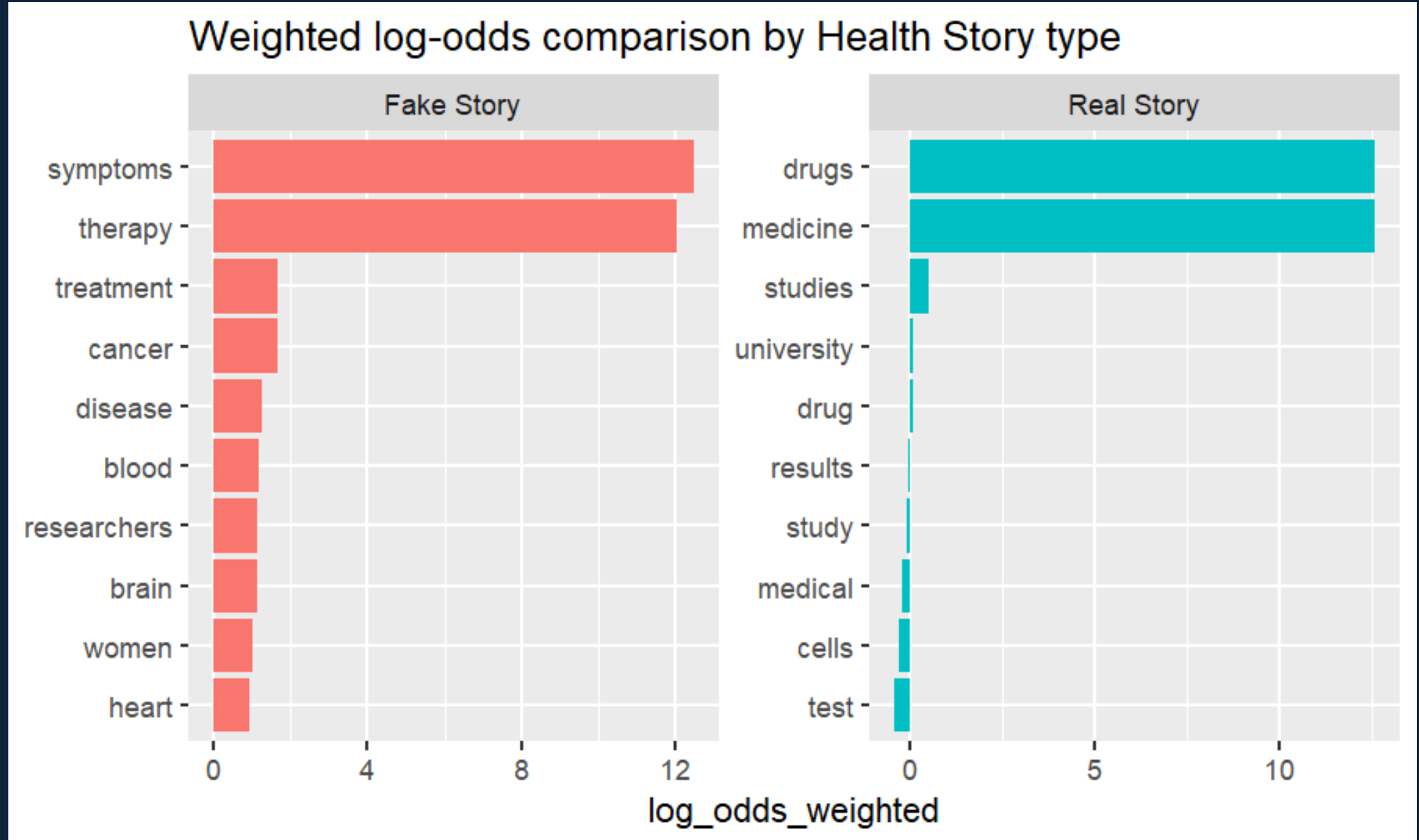


- 1591 documents
- 596 –health release
- 995- health stories

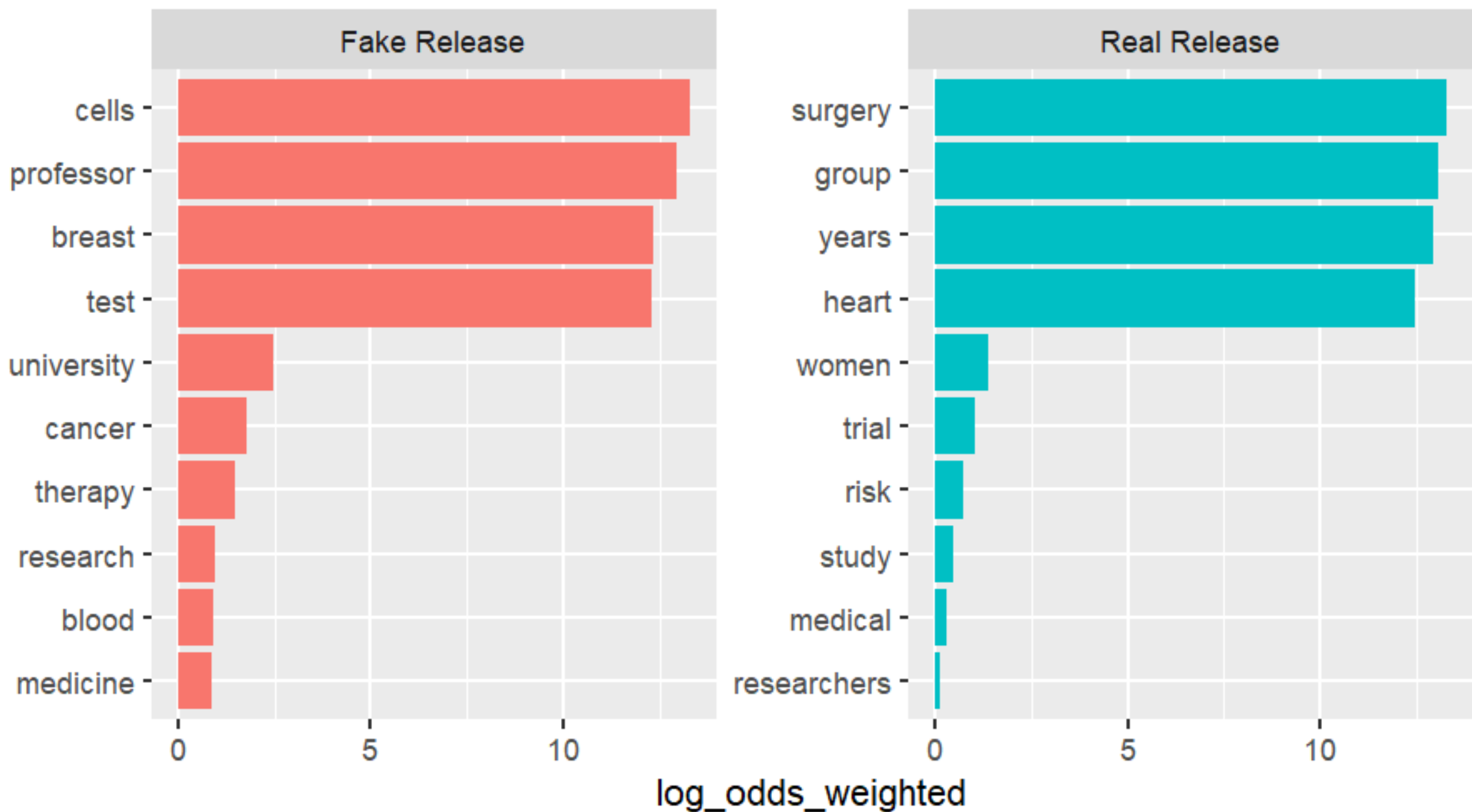
Credibility Criterion

- Less than 3 = fake health information

Methods and Results: Word Ratio Analysis



Weighted log-odds comparison by Health Release type

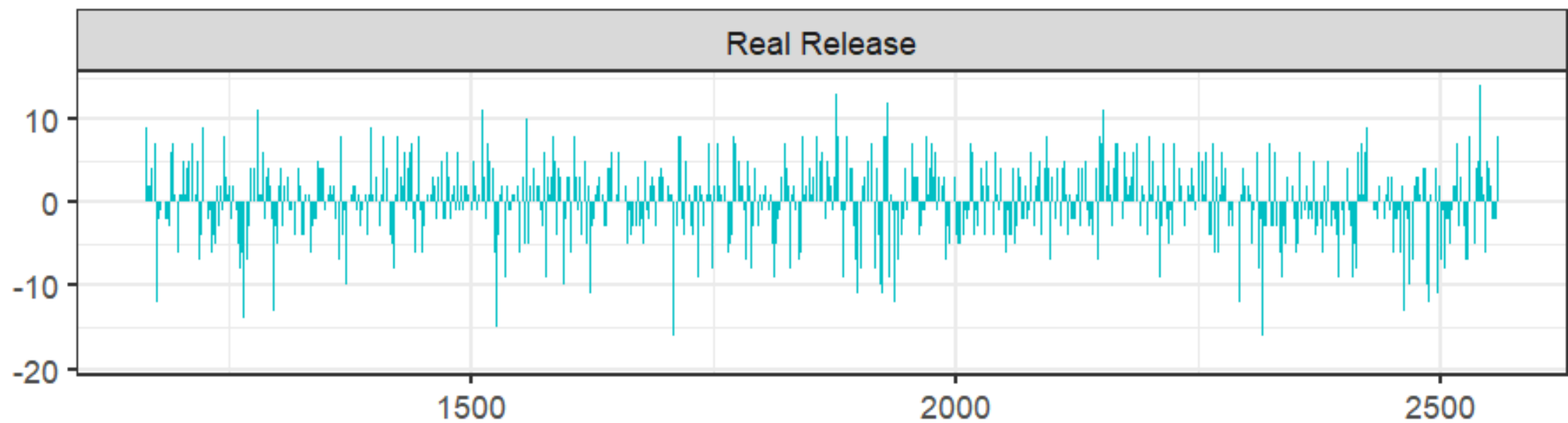
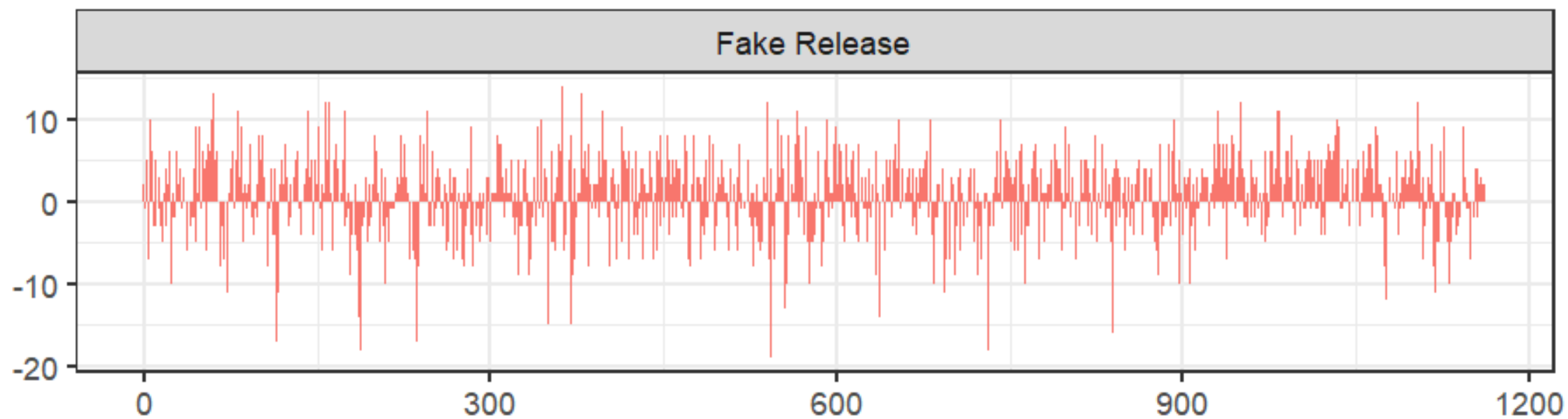


Sentiment Analysis/ Opinion Mining

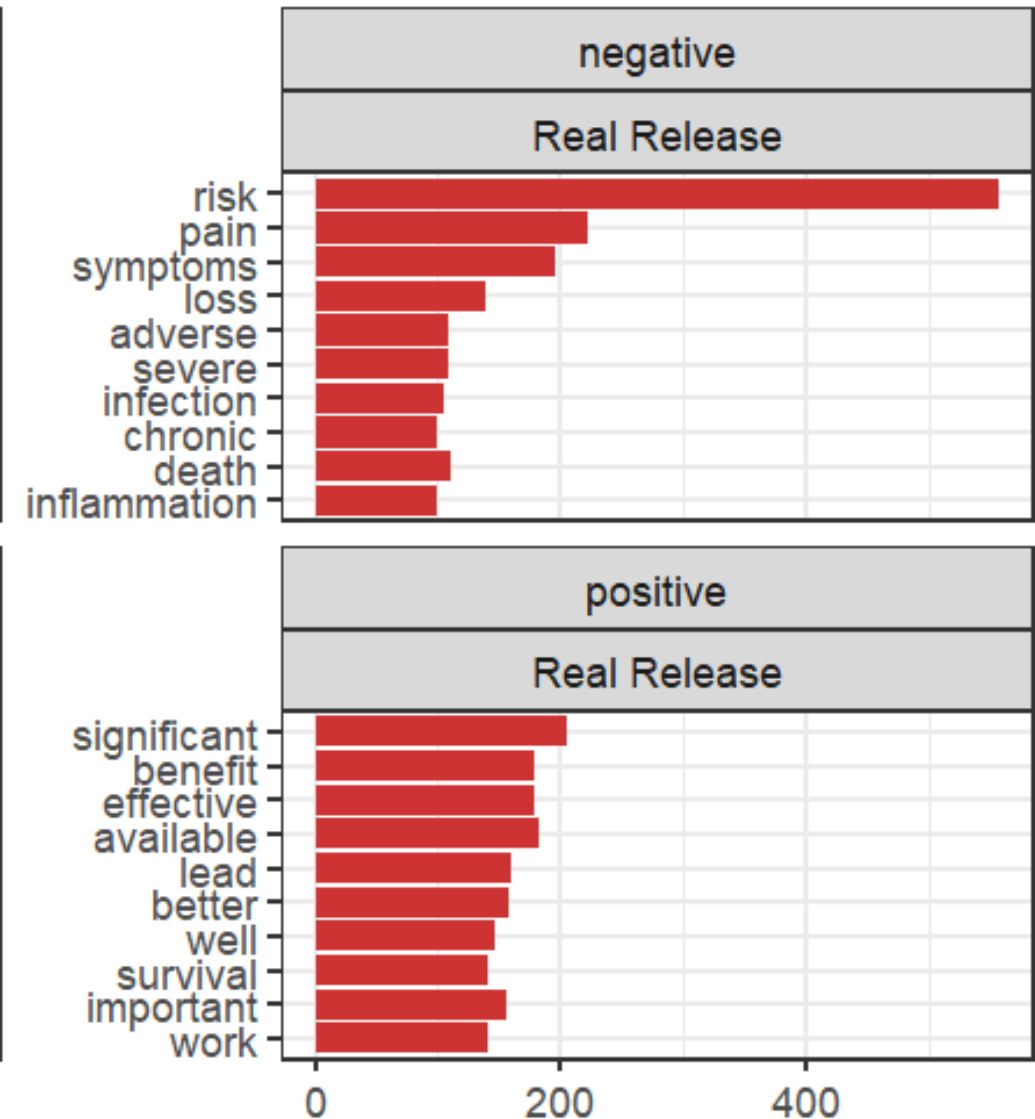
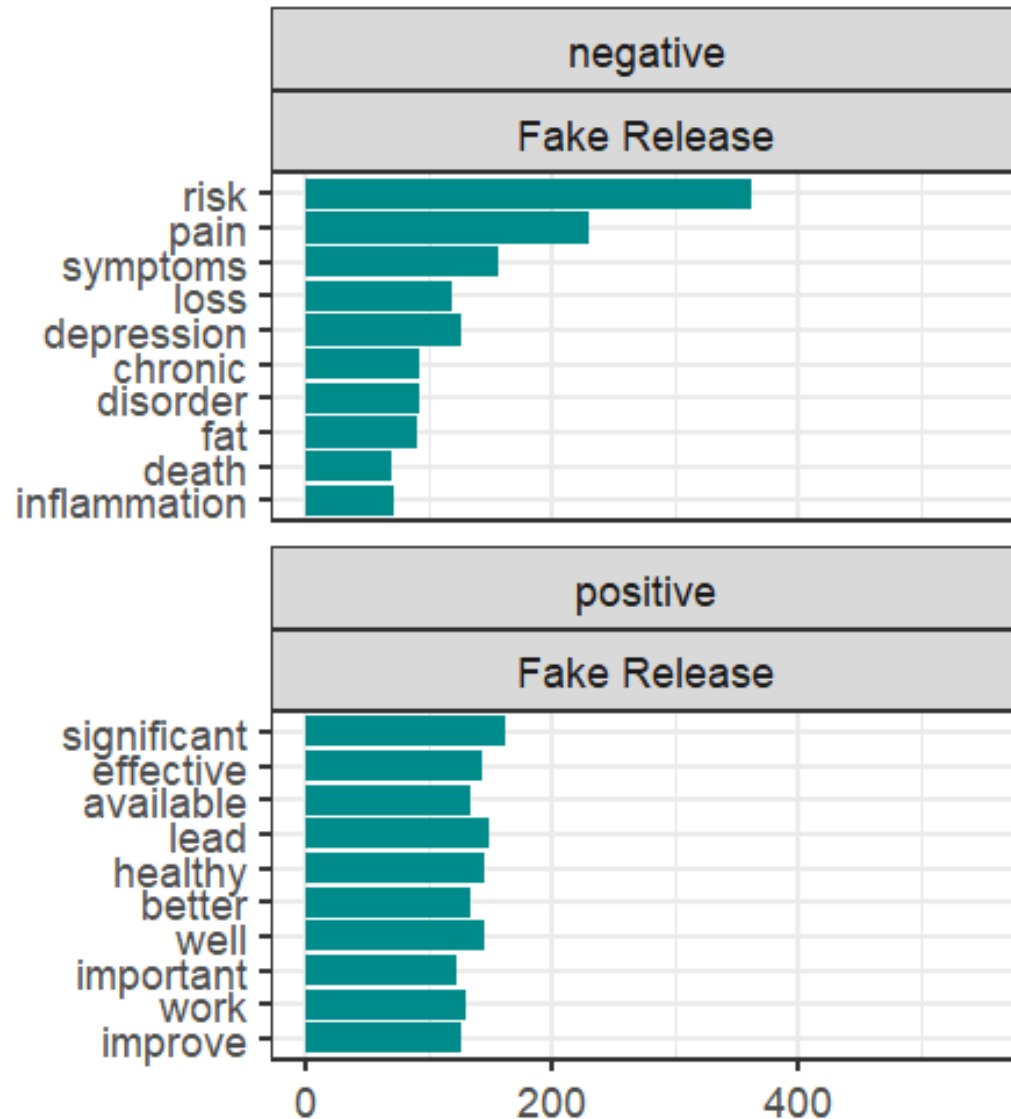
- Sentiment analysis, a form of NLP that evaluates data for neutrality, positivity, or negativity.
- It involves considering individual words and aggregating their sentiments.
- The tidytext package provides access to various sentiment lexicons
- Bing lexicon for emotional intent analysis.
- Bing, structured around single words, assigns negative or positive sentiment to many English words.



Sentiments Score plot by Release Type



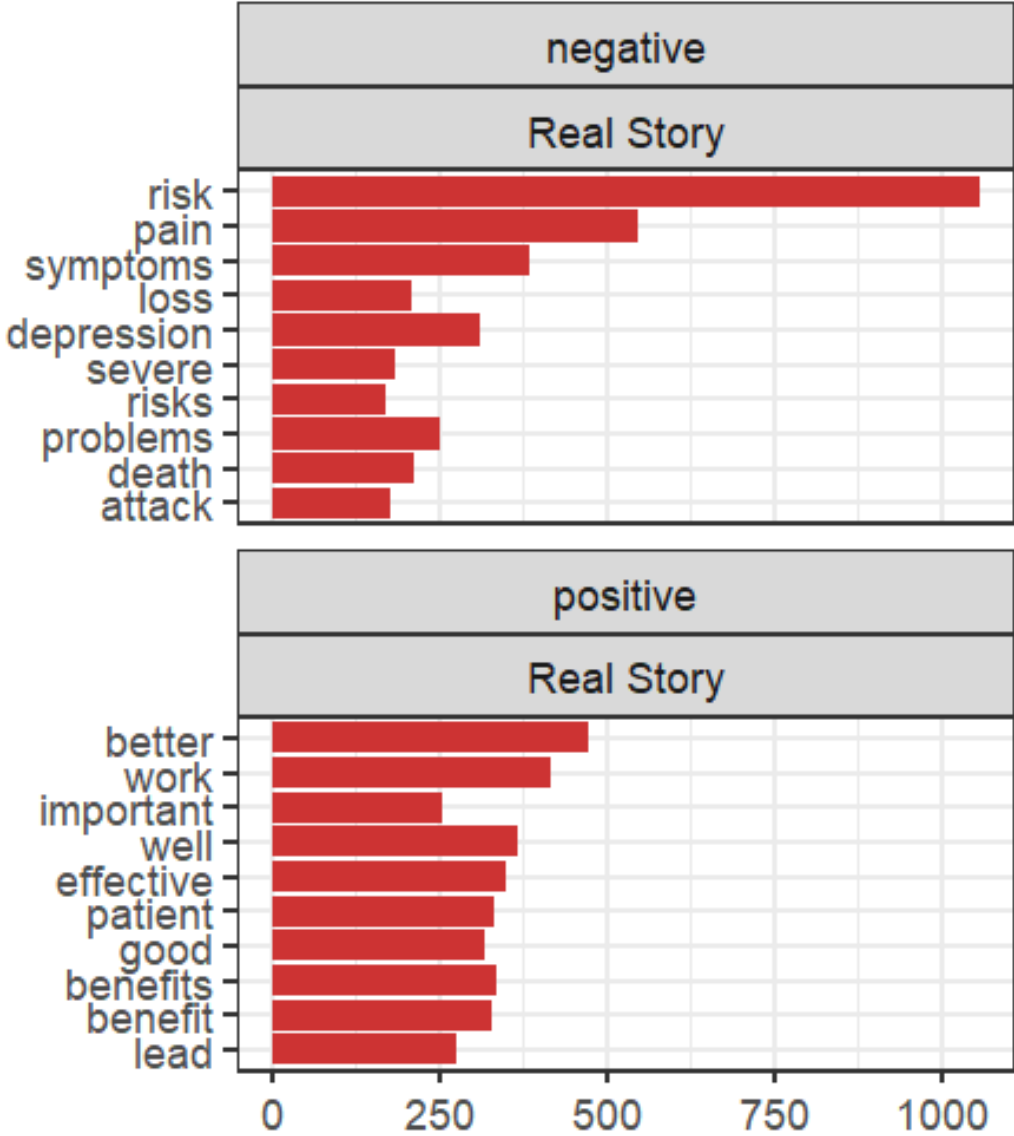
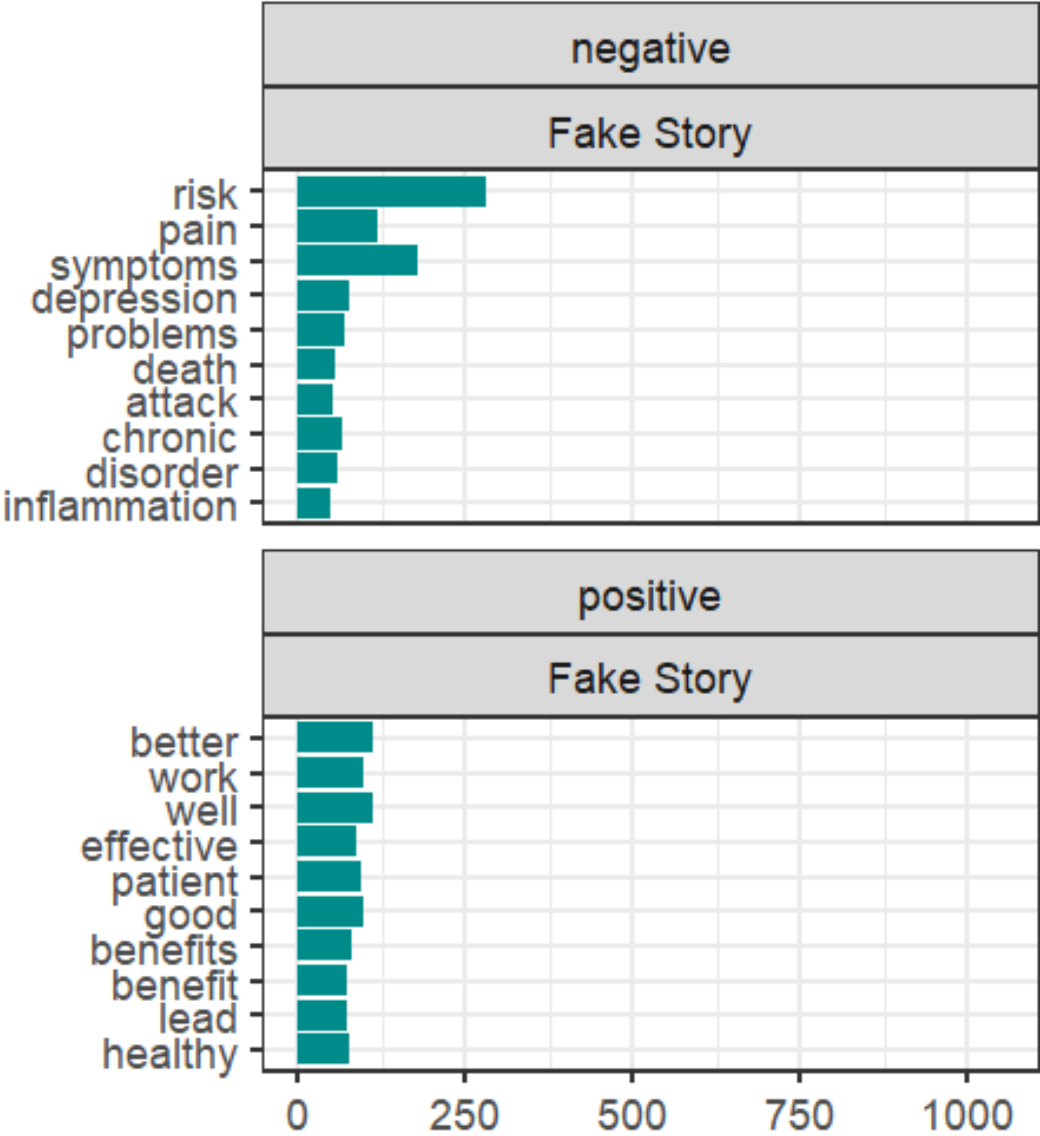
Sentiments Word Analysis by Release Type



Sentiments Score plot by Story Type



Sentiments Word Analysis by Story Type



Topic Modeling



Topic modeling for finding and tracing clusters of words in large bodies of texts.



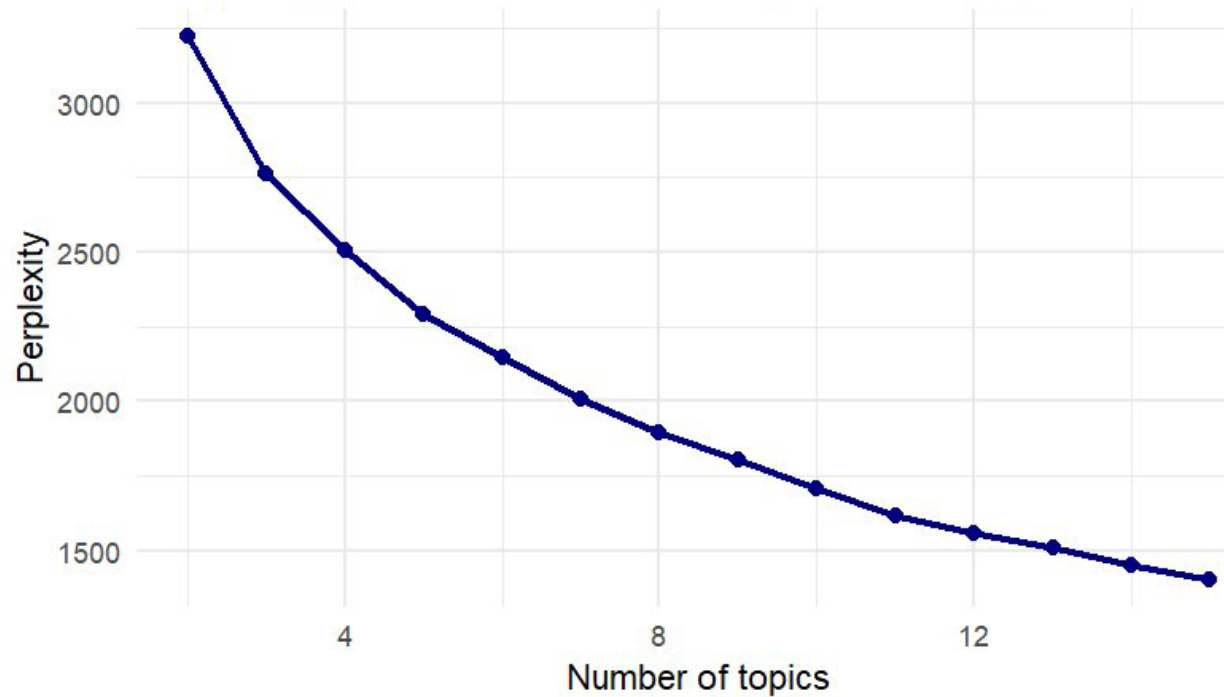
Latent Dirichlet Allocation for its ability to better generalize and its popularity in the humanities fields



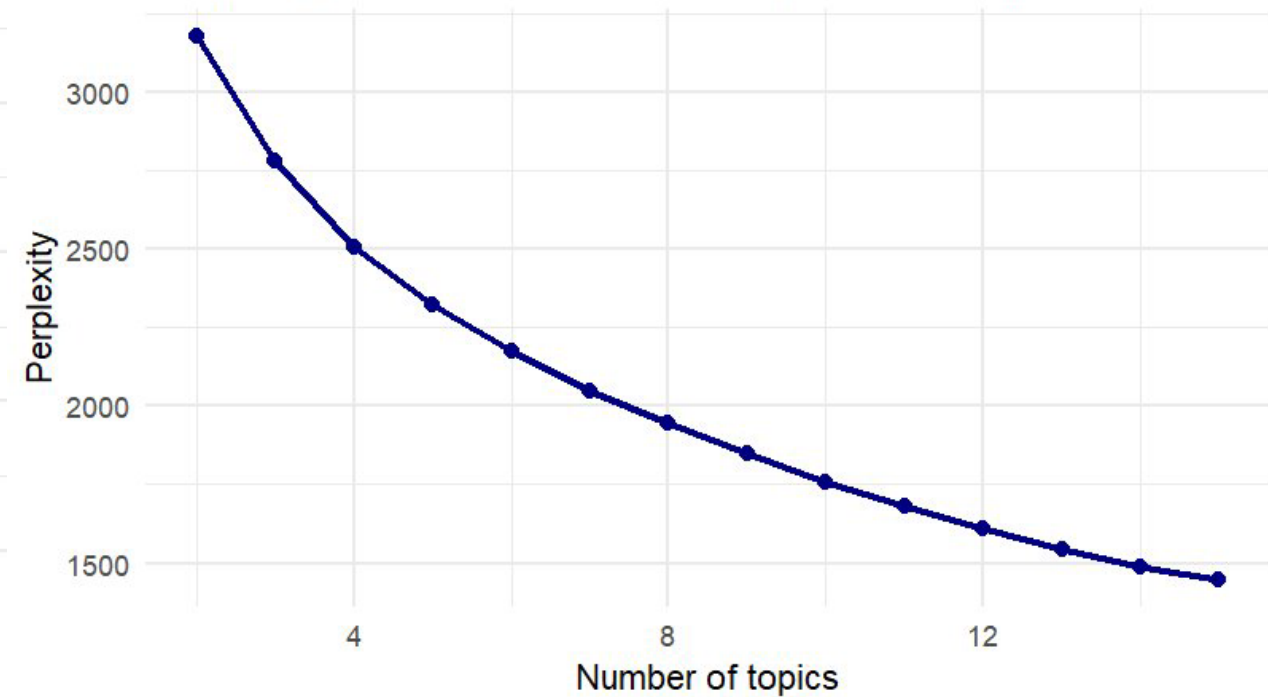
Used elbow plots that mapped the coherence (AKA the perplexity) against the number of topics.

Health Release

Perplexity of different LDA models (Fake Stories)

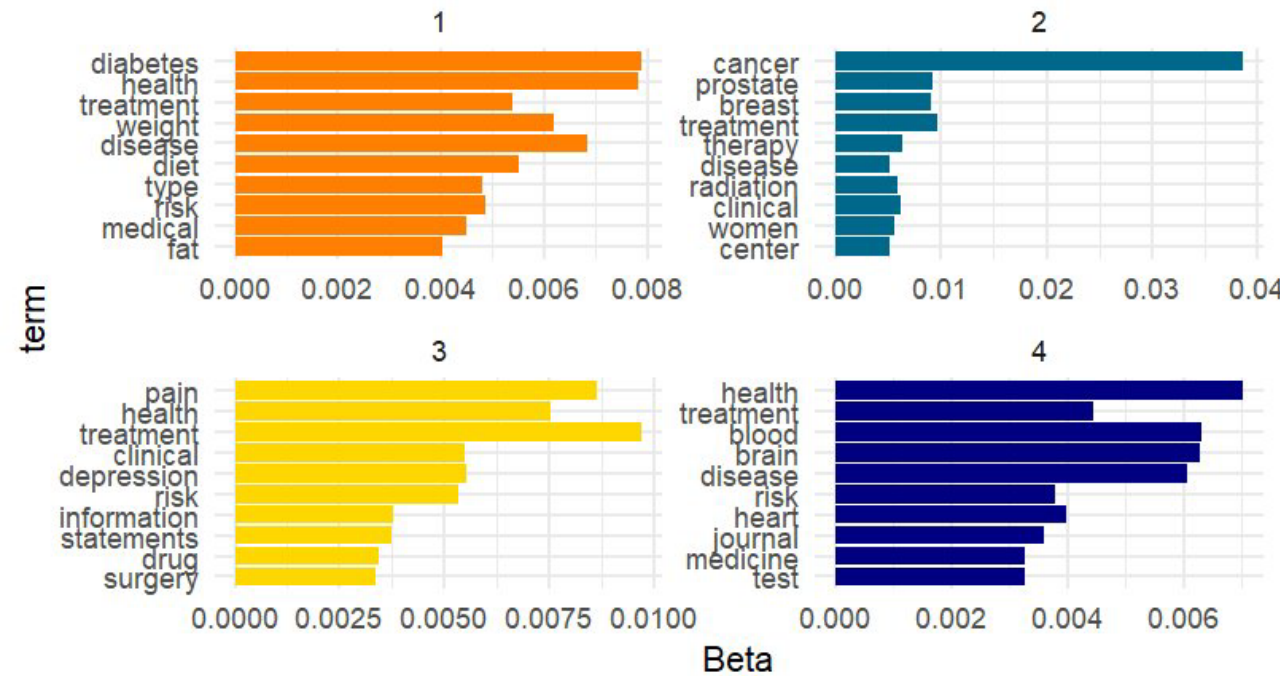


Perplexity of different LDA models (Real Stories)

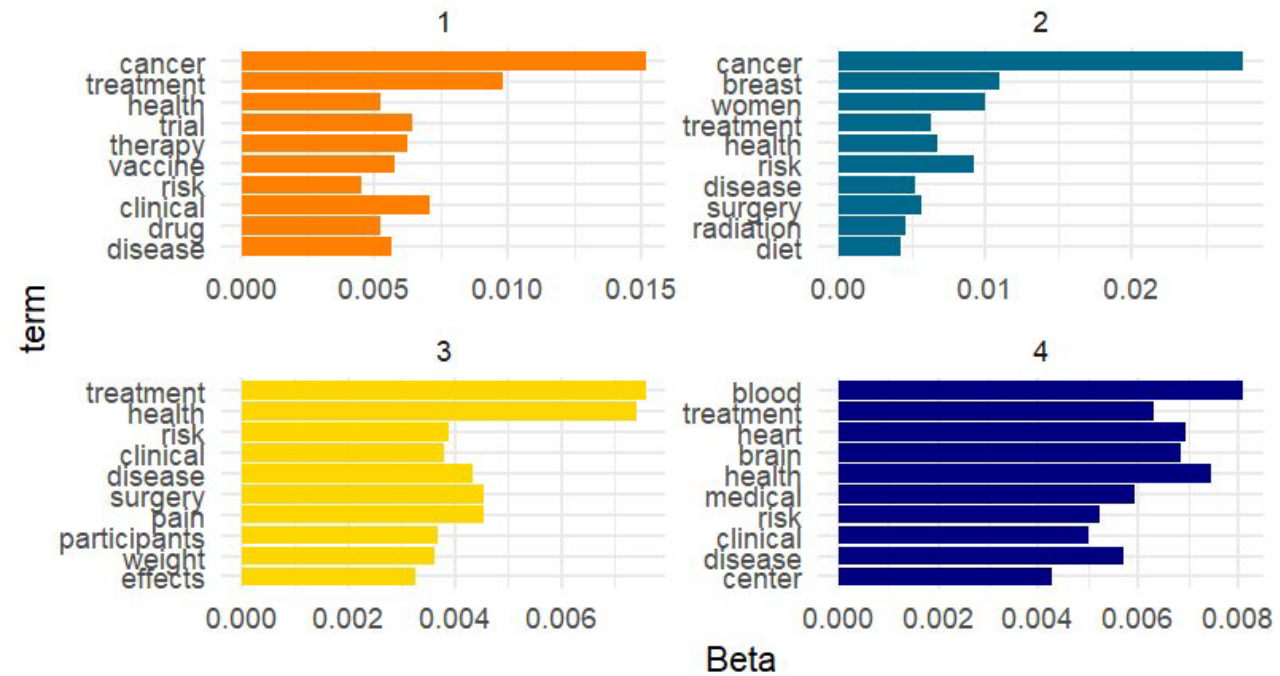


Health Release

Top 10 words per topic (Fake Stories)

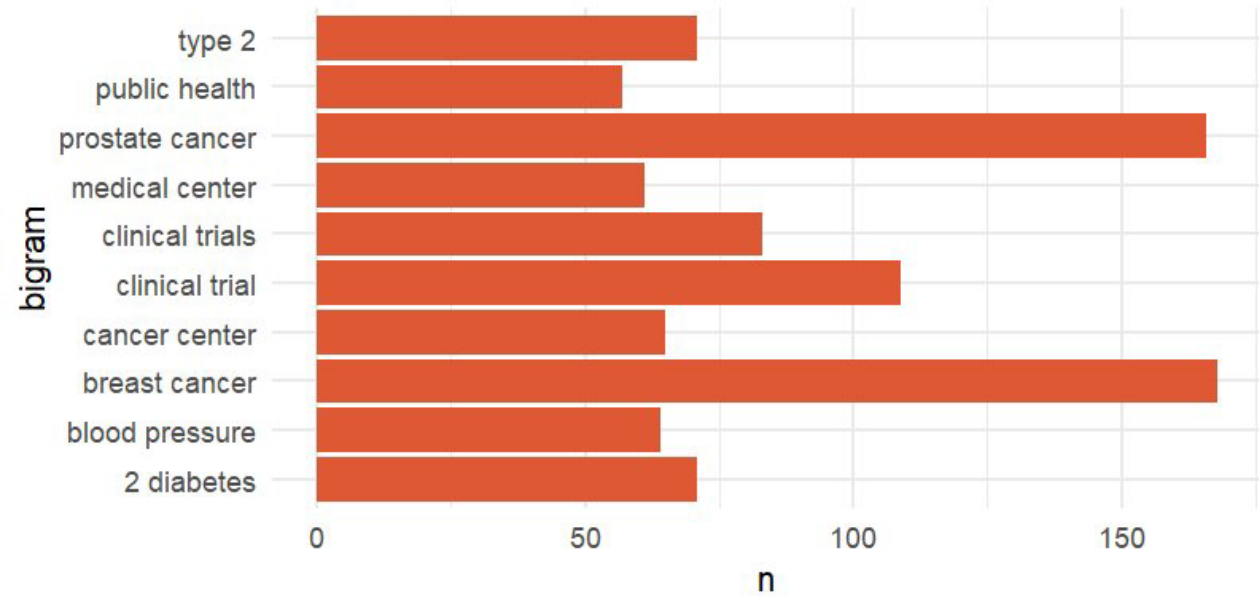


Top 10 words per topic (Real Stories)



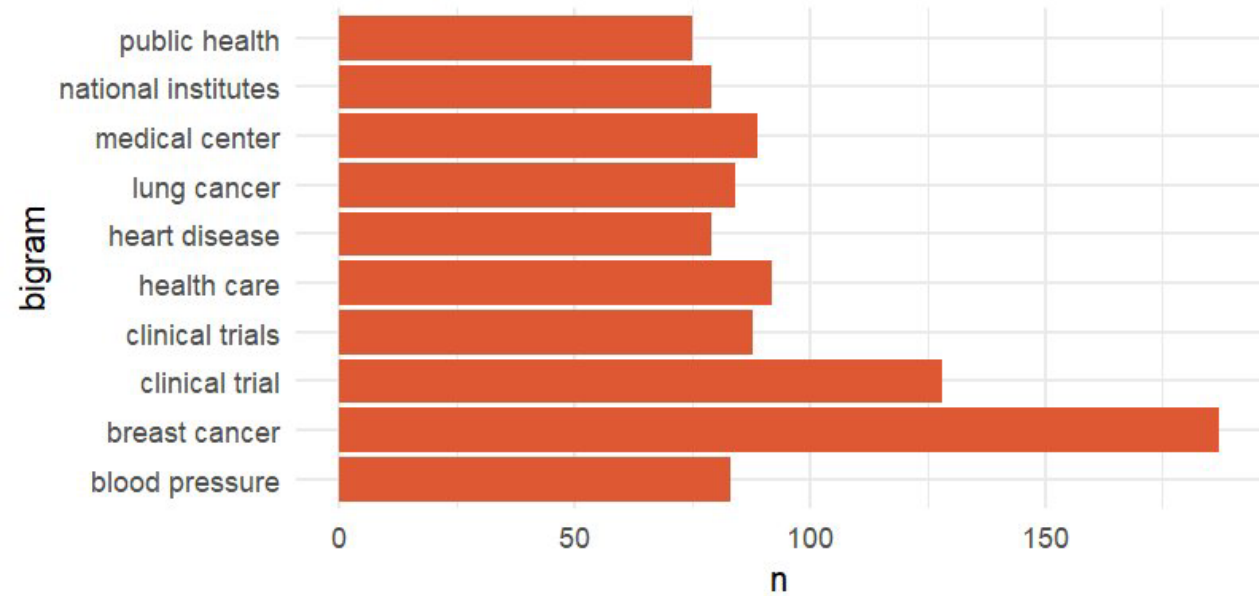
Health Release

Top Bigrams (Fake Stories)
using Tidytext in R



Data Source: Excel File

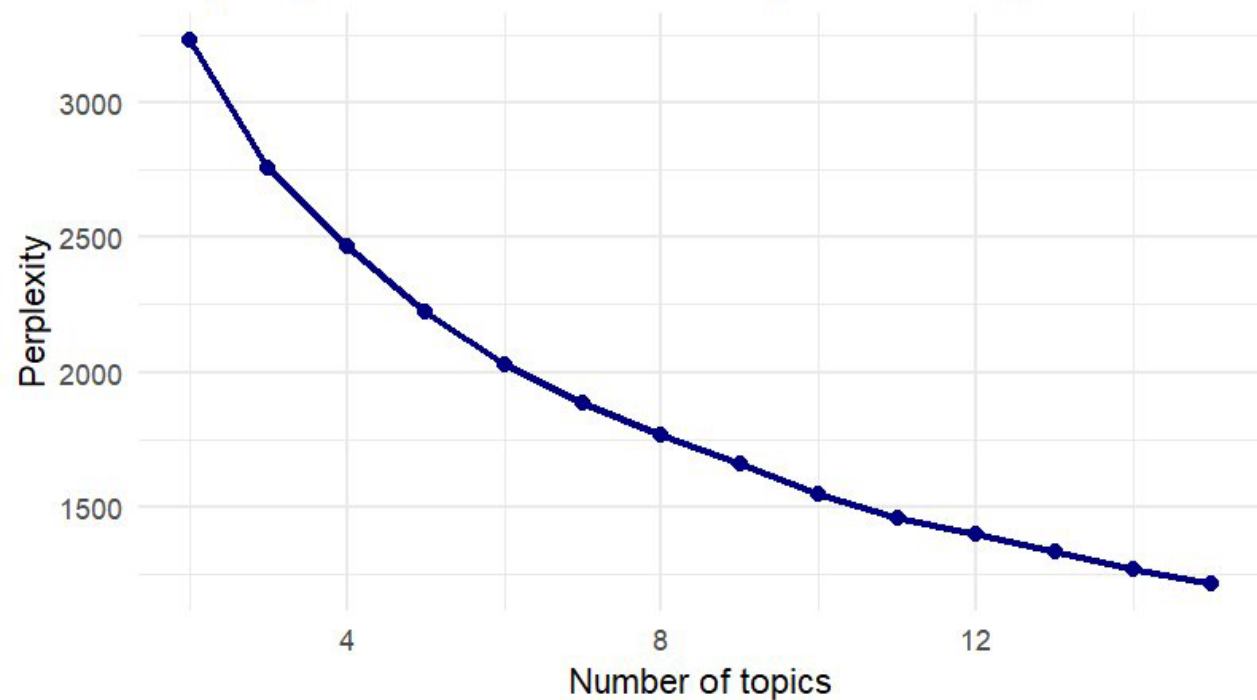
Top Bigrams (Real Stories)
using Tidytext in R



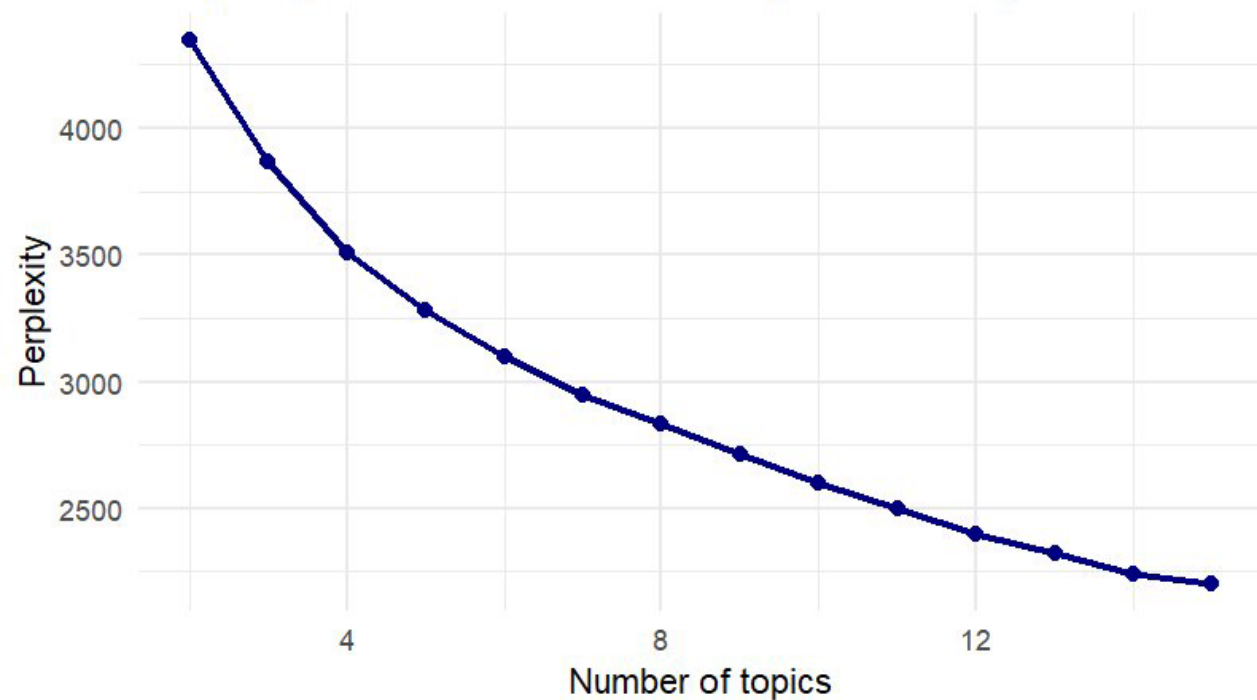
Data Source: Excel File

Health Stories

Perplexity of different LDA models (Fake Stories)

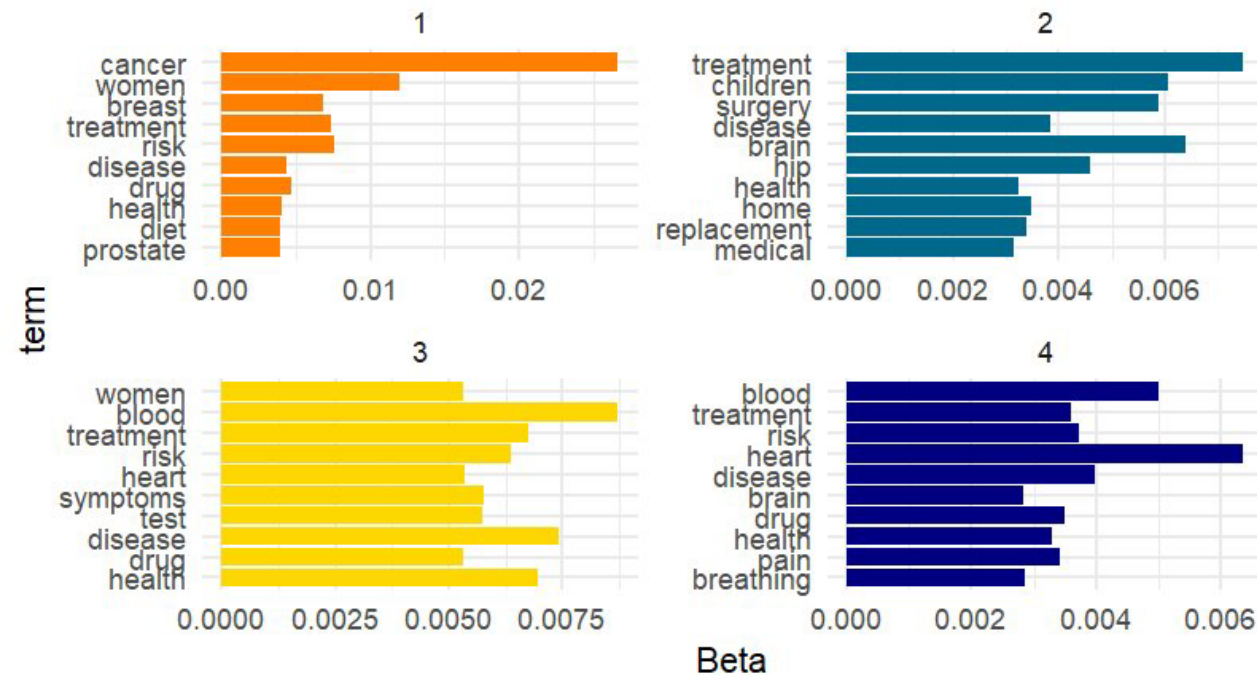


Perplexity of different LDA models (Real Stories)

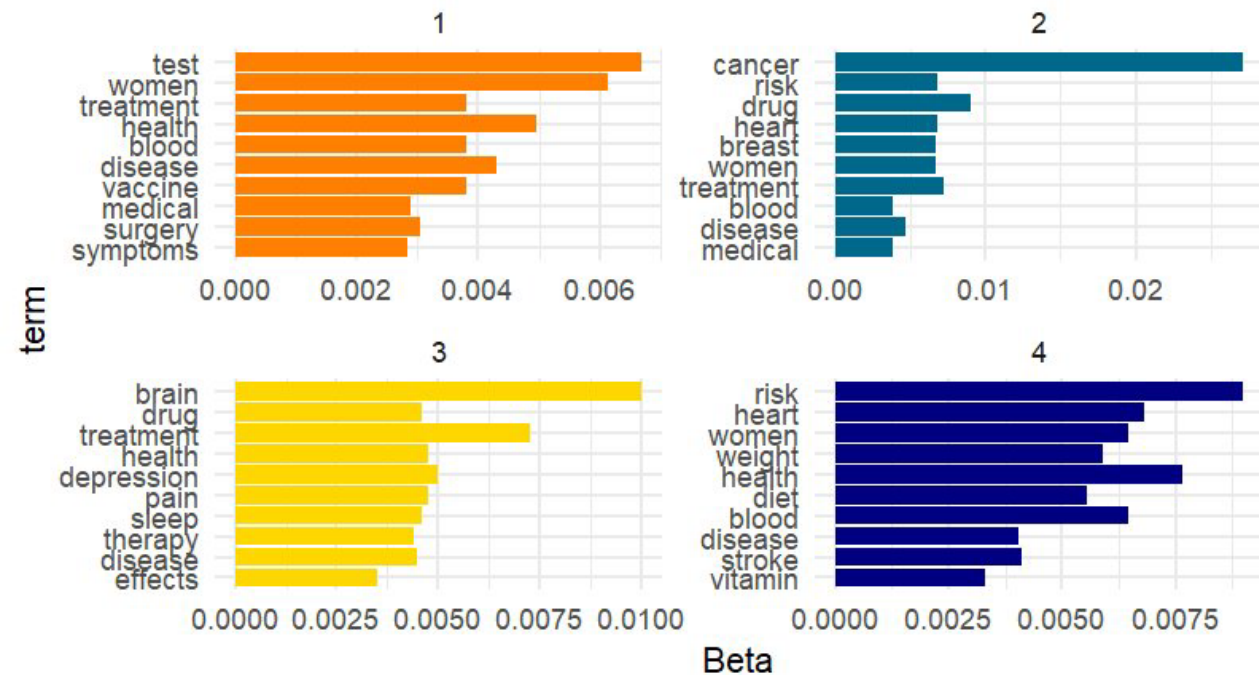


Health Stories

Top 10 words per topic (Fake Stories)

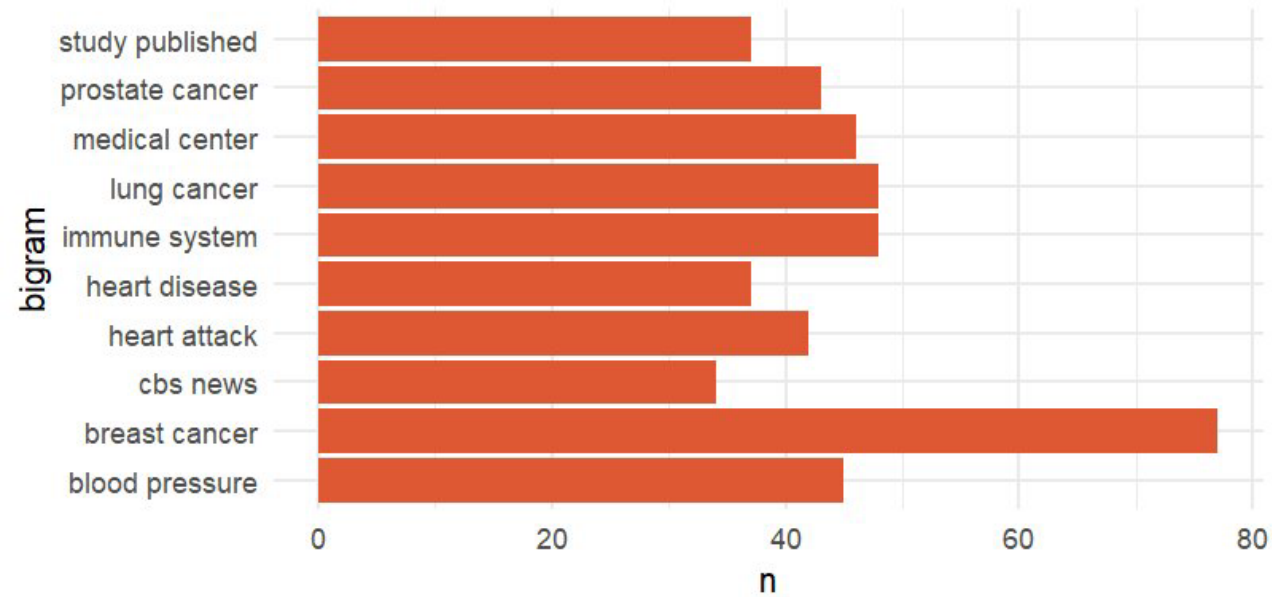


Top 10 words per topic (Real Stories)



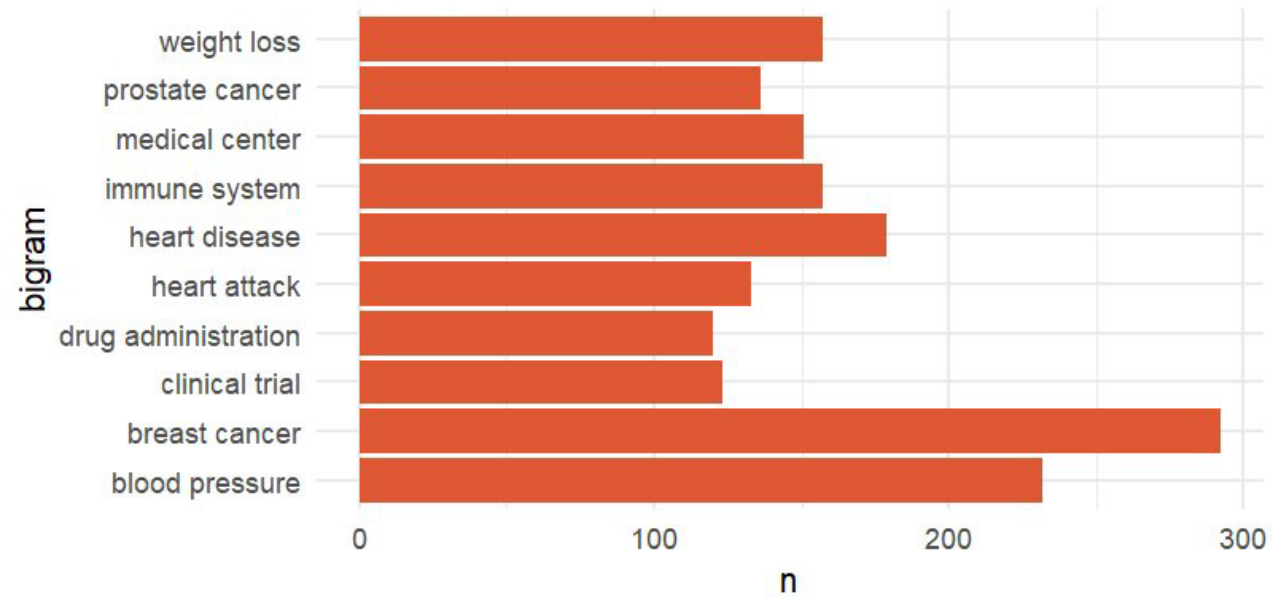
Health Stories

Top Bigrams (Fake Stories)
using Tidytext in R



Data Source: Excel File

Top Bigrams (Real Stories)
using Tidytext in R



Data Source: Excel File

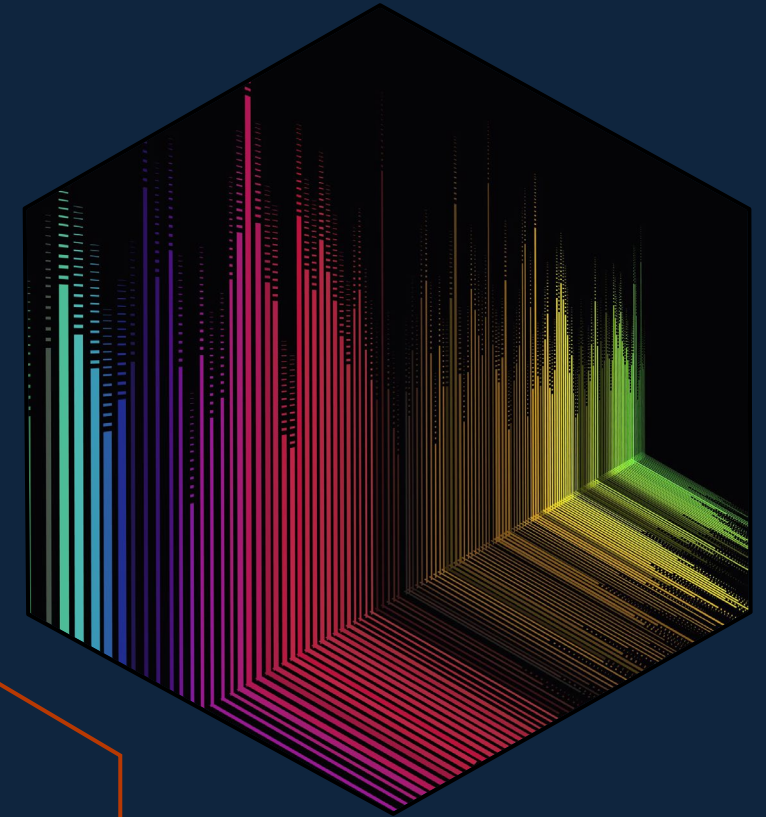
Takeaways and Future Research

1. Mixed support for both our hypotheses

- Some differences in language use between fake and credible health publications.
- Some distinct topics in fake versus real health articles

2. Ideas for future research

- Use a larger sample of documents from variety of sources and conduct a deeper qualitative research
- Compare other sources of health information such as social media



Thank you

