

Disability Trends: An examination of global disability-related terminology

Noah Duggan Erickson

Dept. of Computer Science
Western WA University
Bellingham, WA, USA
email@redacted.com

Brady Deyak

Dept. of Computer Science
Western WA University
Bellingham, WA, USA
email@redacted.com

Raghav Vivek

Dept. of Computer Science
Western WA University
Bellingham, WA, USA
email@redacted.com

Abstract

This document contains the instructions for preparing a manuscript for the proceedings of ACL 2020. The document itself conforms to its own specifications, and is therefore an example of what your manuscript should look like. These instructions should be used for both papers submitted for review and for final versions of accepted papers. Authors are asked to conform to all the directions reported in this document.

1 Introduction

Intro text goes here.

2 Related Work

Related work goes here.

3 Our Approach

Despite the highly multidisciplinary nature of the project, the present work is primarily focused on the natural language processing aspect. The project was divided into three main parts: data acquisition and pre-processing, word vectorization, and sentiment analysis. The data acquisition and pre-processing part of the project involved scraping a large number of news articles from Nexis Uni, converting them to plaintext, and then lemmatizing the text. The word vectorization part of the project involved training a FastText model on the lemmatized text and then visualizing the results. The sentiment analysis part of the project involved training a sentiment analysis model on the lemmatized text and then interpreting the results. The project was conducted using Python and multiple libraries, including Stanza (Qi et al., 2020), Gensim (Řehůřek and Sojka, 2010), TensorFlow (Abadi et al., 2015), and Alibi (Klaise et al., 2021).

4 Experiments

There were multiple experiments run throughout the project. These include a massive scraping operation with subsequently large-scale data pre-processing, a word vectorization model, and an interpretable sentiment analysis model.

4.1 Data Sources

The data for this project was acquired using the export feature of Nexis Uni. Using this feature, a large number of news articles were downloaded as rtf files. These were then converted to plaintext using the `striprtf` (Cyriac, 2023) Python library. Each article's content was then preprocessed using the lemmatization pipeline in the Stanza (Qi et al., 2020) library. This pipeline tokenizes the text, removes stopwords, and lemmatizes the remaining words using a neural seq2seq model. The resulting data was then stored as a list and exported to JSON alongside article metadata.

4.2 Word Vectorization

The word vectorization model used was the FastText model from the Gensim library. After training the model on the preprocessed data from a given decade, the model was used to generate a list of the most similar words to a list of selected words. The results were then visualized using a t-SNE plot. The t-SNE plot was generated using the Plotly library.

4.3 Word Cloud

Another such visual representation of the data after pre-processing was a word cloud generated using the WordCloud library and matplotlib. The word cloud was then generated using word frequencies to set the size of each word.



Figure 1: t-SNE plot of the word vectorization model using the “disab” search term from the 1990s.



Figure 2: t-SNE plot of the word vectorization model using the “disab” search term from the 2020s.

4.4 Sentiment Analysis

BRADY: In this section, discuss the sentiment analysis model. Include the model architecture, training data, and training results. For example, “TensorFlow was used to train a sentiment analysis model similar to the one proposed in this paper. It uses an RNN trained on this dataset, and achieves an average accuracy of $x\%$.”

4.5 Interpretable Sentiment Analysis

RAGHAV: In this section, discuss the interpretations of the sentiment analysis model. For example, how does Alibi work for generating those colors?

5 Results

Multiple visualizations and analyses were conducted on the data. These include a word vectorization model, a document-level sentiment analysis model, and an explanation of said sentiment analysis model.

5.1 Word Vectorization

Word vectorization models were trained on the data for each search term for each decade. This allows for the visualization of not only the most similar words to a given word in that snapshot, but also how the clusters of words shift over time. For example, figures 1 and 2 show the tsne plots for selected words using the “disab” search term from the 1990s and 2020s, respectively.

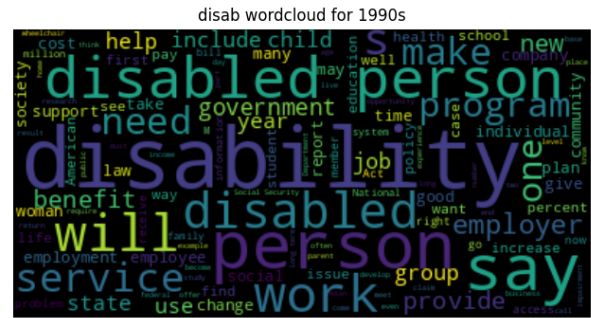


Figure 3: Word cloud of the word frequencies for the “disab” search term from the 1990s.



Figure 4: Word cloud of the word frequencies for the “disab” search term from the 2020s.

5.2 Word Cloud

Word clouds for each term and decade were generated using the word frequencies from the pre-processed data. These word clouds give a visual representation of the most common words in the data. For example, figures 3 and 4 shows the word cloud for the “disab” search term from the 1990s and 2020s.

6 Discussion

As alluded to in our approach, the goal of this project was to create a tool that can be used by other researchers to conduct analyses on disability-related terminology and beyond. As such, the discussion of specific points of interest is outside the scope of this project. However, the results of the analyses shown above represent just the tip of the iceberg of what can be done with this data and these

models. The word vectorization model, for example, can be used to track the evolution of disability-related terminology over time. The sentiment analysis model can be used to track the sentiment of news articles over time. The word cloud can be used to track the most common words in news articles over time. The possibilities are endless.

Furthermore, despite the name of this project, the general idea can be easily extended into other areas simply by using different search terms. For example, one could use the same pipeline to track the evolution of sentiment towards mass transit and transit-oriented development.

7 Conclusion

Due to the short turn-around time of the project, avenues for future work are plentiful. These include a more centralized system for creating the analyses, a more robust sentiment analysis model, and a deeper review of these analyses from a critical disability studies and sociological perspective.

References

- Martín Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro, Greg S. Corrado, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Ian Goodfellow, Andrew Harp, Geoffrey Irving, Michael Isard, Rafal Jozefowicz, Yangqing Jia, Lukasz Kaiser, Manjunath Kudlur, Josh Levenberg, Dan Mané, Mike Schuster, Rajat Monga, Sherry Moore, Derek Murray, Chris Olah, Jonathon Shlens, Benoit Steiner, Ilya Sutskever, Kunal Talwar, Paul Tucker, Vincent Vanhoucke, Vijay Vasudevan, Fernanda Viégas, Oriol Vinyals, Pete Warden, Martin Wattenberg, Martin Wicke, Yuan Yu, and Xiaoqiang Zheng. 2015. [TensorFlow, Large-scale machine learning on heterogeneous systems](#).
- Joshy Cyriac. 2023. `striptrf`: Stripping rtf to plain old text.
- Janis Klaise, Arnaud Van Looveren, Giovanni Vacanti, and Alexandru Coca. 2021. [Alibi Explain: Algorithms for Explaining Machine Learning Models](#). *Journal of Machine Learning Research*, 22(181):1–7.
- Peng Qi, Yuhao Zhang, Yuhui Zhang, Jason Bolton, and Christopher D. Manning. 2020. Stanza: A Python natural language processing toolkit for many human languages. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: System Demonstrations*.
- Radim Řehůřek and Petr Sojka. 2010. [Software Framework for Topic Modelling with Large Corpora](#). Proceedings of LREC 2010 workshop New Challenges for NLP Frameworks, pages 45–50, Valetta, MT. University of Malta.