

Shon Liskov

CS 4395.001

Professor Mazidi

11/9/2022

ACL Summary:

COVID-19 and Misinformation: A Large-Scale Lexical Analysis on Twitter

Dimosthenis Antypas, David Rogers, Alun Preece, and Jose Camacho-Collados

School of Computer Science and Informatics & Crime and Security Research Institute

Cardiff University, United Kingdom

Social media often spreads misinformation. Ever since the pandemic started, people on Twitter tweet out false information about the virus which is harm to public health. Various topics are shared between several twitter users daily. And a good amount of time is spent spreading misinformation about their opinions of COVID-19. Within Twitter there is no system that filters out misinformation so anybody can write their own “facts” without actual prove that it’s true. COVID-19 has only made health related claims worse. The researchers wanted to find a way to distinguish between tweets that are misguided and ones that are generic that doesn’t sound bias. They took a bunch of tweets and compile them into separate corpora which was taken from Twitter data with misinformation related corpus collected from Social Media analysis platform known as Sentinel and generic tweets taken random that is COVID related. An analysis is done to find patterns in respect to vocab usage only using lexical features. The misinformed tweets come out to be more negative and limited in word choice than generic tweets about the virus. Many different models were used to classify between misinformation and generic tweets which includes: SVM, BERT, CNN, Naïve Bayes, etc.

The misinformation corpus was extracted from another study which collects search terms in multiple languages. They can find these types of tweets by finding replies that will say that the tweet is either misinformed, fake news, consist of lies, etc. A table was made to calculate the frequency of features between the two classes. Misinformed tweets tend to be longer, has greater user mentions, and percentage usage of exclamation marks are 62% higher. To find uniqueness between the two classes, TTR/MTLD is used for coming up with a percentage that determines whichever corpora is more diverse than the other. They also used lexical specificity to find the frequency of words throughout three months between both classes (corpora). For the misinformation corpus, the top 3 words were “uncover, deep, and medium”. For the generic corpus, the top 3 words were “confirm, suga, and case”.

Some unique contributions include showing tables of results to certain tasks like when they compare between two classes, they show the accuracies between SVM AND BERT since both have been known to have higher accuracies than any of the other models. Also, a table that shows lexical specificity is done to show frequencies of different words that are classified as misinformation or generic. Lastly, the researchers used a total of twenty-four sources to combine different ideas and bring together to fully express the expected results and share what failed and succeeded during the process of the work with classifying misinformed tweets versus generic tweets.

Dimosthenis Antypas is a teaching associate has done lots of research in data science/machine learning with computational analysis work through Twitter API and has written five papers talking about how he and other authors can solve problems using machine learning models pertaining to Twitter, plus using NLP tasks to figure out lifelong solutions to problems in social media. Jose Camacho-Collados is a senior lecturer with more experience in the computer science with specialty in NLP and has written a book about recent trends in distributional semantics and NLP. He has over 90 citations from NLP to AI, etc. Alun Preece is a deputy head of the computer science department at the university and has over

400 citations. David Rogers has over 263 citations. All these authors have in common is NLP and working with NLP tasks for variety of different problems in social media or the internet in general.