

# Medical Relevancy of Cancer-Related Tweets and Its Relation to Misinformation

Melanie McCord, Fahmida Hamid

New College of Florida  
Division of Natural Sciences  
{melanie.mccord19, fhamid}@ncf.edu

## Abstract

Social media is one of the most dominant ways of spreading information. Still, unfortunately, these open platforms provide ways to spreading misinformation which can be extremely dangerous, especially when relevant to sensitive issues such as health-related information. Hence such platforms require an effective autonomous misinformation detection mechanism. Understanding the data is one of the necessary artifacts for building such a mechanism. In this work, we attempted to determine the medical relevancy of cancer-related tweets and explore whether they contain misinformation. We created a dataset of roughly 500 tweets and labeled them according to their medical relevance: medically relevant, not medically relevant, or unrelated to cancer. We ran logistic regression and support vector machine models on them. The highest proportion of correctly identified “medically relevant” tweets, i.e., accuracy, was 0.795. Our analysis hints at some features and factors that can automatically improve cancer-relevant and non-relevant tweet detection.

## Introduction

Misinformation is prevalent on social media. Users share fake news stories that tend to propagate more than accurate news, mainly because the fake news often elicits negative emotions, such as fear or anger (Shu et al. 2017). When exposed to misinformation, people are often influenced by it (Lewandowsky et al. 2012). Healthcare misinformation and particularly cancer misinformation can have deadly consequences, such as leading people to use dangerous and ineffective remedies, delay effective treatments, or use ineffective strategies to treat cancer. This project aimed to identify medical relevance in tweets among cancer-related tweets and relate it to the detection of misinformation among tweets.

Previous work on detecting cancer misinformation includes using machine learning models with the DETERRENT (Cui et al. 2020) cancer dataset and the Extent of Misinformation (Bal et al. 2020) dataset. Our dataset focuses on tweet engagements and compares the fake news claims among medically relevant versus non-medically relevant tweets.

Copyright © 2023 by the authors. All rights reserved.

## Experimental Design

### Dataset Construction

The seed keywords we used for tweet collection are “prevent”, “cure”, and “cancer”. We then used a collection of keywords gathered from medical journals (an automatically scraped dataset constructed with PubMed papers) to expand the search. Additionally, we used the noun phrases in the DETERRENT news articles to search for tweet-related misinformation. In total, we identified 42,140 tweets. Of these 42,140 tweets, we randomly selected 2000 for initial labeling, of which roughly 500 have been labeled (Table 1).

| Tweet Category         | Count |
|------------------------|-------|
| Medically relevant     | 201   |
| Not medically relevant | 182   |
| Unrelated to disease   | 111   |

Table 1: Tweet categories by count

### Initial Labeling: Medical Relevancy

Initially tweets were placed into three categories: medically relevant, not medically relevant, and unrelated to disease. Medically relevant tweets include tweets that make claims about cancer, share cancer-related studies, or address medical information. Not medically relevant tweets discuss issues, such as cancer instances in people or describe political figures. Unrelated tweets include tweets that are not related to the disease cancer, such as those that are related to zodiac signs or politics. Table 2 shows some instances.

### Secondary Labeling: Sub-Categories

Each tweet was labeled according to whether it contained misinformation or not. To identify misinformation, the tweets were fact-checked using PubMed articles, cancer.gov, as well as reputable news sources. Generally, tweets that were considered medically relevant were either true or false, with most of the non-medically relevant or unrelated tweets being unverifiable.

Subsequently, tweets were further classified into categories depending on the content, motivation, and finally, whether they contained misinformation. For the scope of this

| Category             | Sample Tweet   |
|----------------------|--|
| Medically relevant   | Among women, gynecologic cancers that affect the reproductive system are one of the most common. It's important to be aware of the risks and symptoms, so you can see guidance and treatment options from a medical professional. #gynecologiccancers #Cancers   |
| Medically irrelevant | Lovely evening with @providence donors, patients, scientists physicians in support of cancer research @ChilesResearch. Thank you to our gracious hosts, Aryes wine-maker Brad McElroy his wife Kathleen, for their hospitality amazing vintage wine-PROVIDENCIA! <a href="https://t.co/bIMWZbmzJ4">https://t.co/bIMWZbmzJ4</a> |
| Irrelevant           | @OccupyDemocrats My 50 year republican voter viewpoint: Many republicans put party before country, some are quiet, hope it all goes away, and some that know we have to cut out the cancer in our party. We are not the party of Eisenhower, Reagan, the Bushes, or McCain. We can't learn from Cheney.                        |

Table 2: Examples of Tweets in Each Category

study, we decided to filter them down into 6 categories: Personal, Claims, Advertisement, Political, Articles, and Miscellaneous. Table 3 and figure 1 describes categorical relevance.

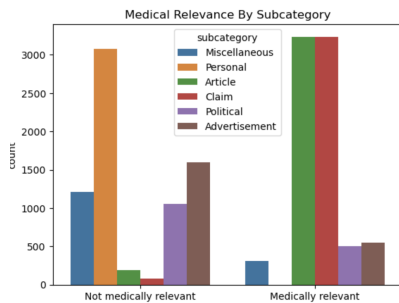


Figure 1: Medical Relevancy of the Sub Categories

### Labeling Tweets with Misinformation

One of the significant difficulties in labeling the tweets was the vast subcategories of misinformation. While some tweets are blatantly untrue, there are also several ways a tweet can be somewhat actual but contain misinformation or unverifiable information. For example, a tweet may be partially true, such as claiming that birth control pills increase cancer risk. While this is true for breast cancer, it also decreases the risk

| Sub-Category  | Sample Tweet   |
|---------------|--|
| Personal      | Both Austin & his Shannara character 'Wil' mothers dies from cancer. A tough script.   |
| Claims        | Deficiency of Calcium & Vitamin D leads Cancer and Depression.   |
| Advertisement | Want to find out how to make Cancer-Fighting Lifestyle Recipe for Zucchini Bites? It's as easy as "clicking" the link below. #Thanks to @StanfordMed!  |
| Political     | If O'Connor takes this step—which he gestures toward in today's ruling—private health insurers could refuse to cover huge range of preventive care, including vaccines, cancer screenings, STI tests, pregnancy care ... the list goes WAY beyond PrEP coverage. <a href="https://t.co/BChDAE3Yhr">https://t.co/BChDAE3Yhr</a> |
| Articles      | #US lab at the centre of legal fight over #Zantac and #cancer <a href="https://t.co/7WYpDeSkOe">https://t.co/7WYpDeSkOe</a>  |
| Miscellaneous | Do yourselves a favor and look at \$PPCB Dropped a PR this morning stating success about cancer treatment and reversing tumor development. Bio-med and healthcare play.  |

Table 3: Examples of the Annotations

| Label        | Tweet Count |
|--------------|-------------|
| True         | 140         |
| False        | 44          |
| Can't verify | 116         |

Table 4: Labeled verified tweets

of other cancers. Additionally, scientific works can be inconclusive. Due to the dataset's limited size, we have labeled all misinformation as false. Table 4 shows the total tweet counts for each type (True, False, Unverified).

### Analysis and Discussion

We fitted two supervised models: support vector machine and logistic regression with a TF-IDF transformer through sci-kit learn. The support vector machine performed the best. Table 5 shows our results. Overall, this study represents an important classification scheme to categorize false or misleading tweets.

| Model               | Accuracy | Recall | Precision | F1-Score |
|---------------------|----------|--------|-----------|----------|
| SVM                 | 0.785    | 0.804  | 0.755     | 0.779    |
| Logistic Regression | 0.795    | 0.717  | 0.825     | 0.767    |

Table 5: Model results

## References

- Bal, R.; Sinha, S.; Dutta, S.; Joshi, R.; Ghosh, S.; and Dutt, R. 2020. Analysing the Extent of Misinformation in Cancer Related Tweets. *Proceedings of the International AAAI Conference on Web and Social Media* 14:924–928.
- Cui, L.; Seo, H.; Tabar, M.; Ma, F.; Wang, S.; and Lee, D. 2020. DETERRENT: Knowledge Guided Graph Attention Network for Detecting Healthcare Misinformation. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, KDD '20*, 492–502. New York, NY, USA: Association for Computing Machinery.
- Lewandowsky, S.; Ecker, U. K. H.; Seifert, C. M.; Schwarz, N.; and Cook, J. 2012. Misinformation and its correction: Continued influence and successful debiasing. *Psychological Science in the Public Interest* 13(3):106–131. PMID: 26173286.
- Shu, K.; Sliva, A.; Wang, S.; Tang, J.; and Liu, H. 2017. Fake News Detection on Social Media: A Data Mining Perspective. *ACM SIGKDD Explorations Newsletter* 19(1):22–36.