

DETECTING HEALTH MISINFORMATION IN WEB PAGE TEXT USING DEEP LEARNING METHODS

Dione Morales

Bachelor of Engineering
Computer Engineering Stream



School of Engineering
Macquarie University

November XX, 2018

Supervisor: Associate Professor Adam Dunn

ACKNOWLEDGMENTS

I would like to acknowledge ...

STATEMENT OF CANDIDATE

I, (insert name here), declare that this report, submitted as part of the requirement for the award of Bachelor of Engineering in the School of Engineering, Macquarie University, is entirely my own work unless otherwise referenced or acknowledged. This document has not been submitted for qualification or assessment at any academic institution.

Student's Name:

Student's Signature:

Date:

ABSTRACT

This is where you write your abstract ...

Contents

Acknowledgments	iii
Abstract	vii
Table of Contents	ix
List of Figures	xi
List of Tables	xiii
1 Introduction	1
1.1 Project Overview	1
1.2 Motivation	2
1.3 Aims	2
2 Background and Related Work	5
2.1 Assessing the Quality and Credibility of Health Information	5
2.1.1 Checklists and Criteria	5
2.1.2 Studies Applying Checklists to Health Information Online	6
2.2 Document Classification Methods	6
2.2.1 Representing Text as Features	7
2.2.2 Machine Learning Methods for Document Classification	7
2.2.3 Deep Learning Methods for Document Classification	10
2.2.4 Deep Learning Models	10
2.2.5 Classifying Documents with Limited Labelled Examples	10
2.3 Discussion	11
2.3.1 Scaling the Assessment of Quality and Credibility of Health Information	11
2.3.2 Using Deep Learning Methods to Automate Quality and Credibility Assessments	11
3 Proposed Approach	13
3.1 Rationale	13
3.2 Credibility Criteria	13

3.3	Study Data	13
3.4	System Model	13
3.5	Experiments	13
3.6	Outcome Measures	13
4	Conclusions and Future Work	15
4.1	Conclusions	15
4.2	Future Work	15
5	Abbreviations	17
A	name of appendix A	19
A.1	Overview	19
A.2	Name of this section	19
B	name of appendix B	21
B.1	Overview	21
B.2	Name of this section	21
	Bibliography	21

List of Figures

2.1	What is the Best Separating Hyperplane? [4]	9
2.2	(Left) Non-linearly distributed data of two classes in a 2D representation space. (Right) Linearly separable data of the same classes in a 3D representation space after the application of a kernel [3].	9

List of Tables

Chapter 1

Introduction

Things to talk about:

- focus on the scenario if we were able to evaluate the quality/credibility of information at scale
- this would allow us to identify information (and the communities that are susceptible to it) - allowing us to intervene and minimise the spread and effect on behaviour
- being able to do this at scales either requires a lot of experts assessing the online health information (resource intensive) or have an algorithm that can learn to do it automatically (less resource intensive)

With the popularity and ubiquity of social media platforms in today's society, the amount and rate at which information propagates online greatly outnumbers the resources available that can evaluate the quality and credibility of the information that gets shared. This is becoming an increasingly growing issue due to the clickbait model that is commonly adopted by social media platforms and the lack of rigour surrounding the publishing of content online [16], causing an increase in the number of articles that contain misinformed content [17] [18]. The lack of resources that attempt to minimise the spread of misinformation can be attributed to the expert-level knowledge required to determine the credibility of information since a variety of factors such as the actual information, information sources, conflicts of interest and writing style of the article must be evaluated to be able to determine its quality and credibility. This issue is further exacerbated in specific domains such as for health-related text, as the spread of misinformation can have a detrimental effect on communities that don't have access to these kinds of resources. Thus, by developing a method that is capable of automatically assessing the credibility of an article in order to

1.1 Project Overview

This section details the scope of the project and its associated outcomes outlining the various tasks that must be accomplished to successfully complete the project.

One of the key components required to minimize the propagation of misinformation online is to have the ability of automatically evaluating and quantifying the credibility of articles. However, traditional automated methods - such as machine learning-based techniques, still require the domain knowledge of experts to be able to develop the features required by the model. Thus, this project aims to investigate the performance of Deep Learning-based (DL) techniques in evaluating the credibility of information within domain-specific articles via the classification of set criteria that have deemed to be highly correlated with articles that have low credibility. Specifically, this project will focus on evaluating the credibility of online health articles related to vaccination due to the commonly misinformed and controversial views associated with its effects [8].

1.2 Motivation

Try to answer some of the following questions:

- Why is misinformation and the spread of low credibility information a problem?
- What is vaccine hesitancy, how many people believe that vaccines are harmful?
Explain specific examples

1.3 Aims

With the primary objective of this project being the evaluation on the effectiveness of deep learning models in determining the credibility of online health-related articles. Due to the complexity of this project, a set of activities - divided into main goals and stretch goals, have been defined to ensure that the completion of this project remains feasible in the given time frame. The completion of all activities categorized as main goals will signal the realization of the primary objective and the completion of the project. Stretch goals are activities of interest that have been identified as non-essential to the completion of the primary objective but (talk about the overarching goal that all stretch goals have in common e.g. understand the model, utilize the model etc.) and will be worked on after the completion of the project.

Main Goals

- Implement and evaluate the performance of machine learning methods for assessing the quality and credibility of vaccine-related text.
- Implement a deep neural network method to assess the quality and credibility of vaccine-related text and compare the performance with the previous approaches.

Stretch Goals

- Evaluate the effect of transfer learning methods on the training time and performance of the proposed deep neural network method.

Chapter 2

Background and Related Work

A literature review has been conducted to develop an understanding on the research that has been done in the assessment of the credibility of information, specifically in the context of information related to health and the limitations and capabilities of machine learning techniques and how it differs from deep learning-based methods for the task of document classification.

2.1 Assessing the Quality and Credibility of Health Information

To have the capability of automating the process of evaluating the quality or credibility of online information, a definition that outlines of what is required by an article to be considered as a credible source of information must be developed. While there has been a significant amount of research that has been done on the development of tools and frameworks that aim to assess the credibility of online health information, there is currently no standardized method or benchmark that is universally used. Tools and frameworks that have been identified to be applicable in the context of this project are: DISCERN [1], HealthNewsReview [2] and Quality Index for health-related Media Reports (QIMR) [19].

2.1.1 Checklists and Criteria

DISCERN

DISCERN is a questionnaire designed to assess the reliability of a publication, it consists of 16 questions each with a Likert scale, ranging from 1 (no) to 5 (yes) and is divided into 3 sections. The first section (questions 1 - 8) investigate the reliability of the information and is comprised of questions such as "Are the aims clear?", "Is it balanced and unbiased?" and "Does it refer to areas of uncertainty?". The second section (questions 9 - 15) assesses the quality of information provided by the publication for treatment choices and is composed of questions such as "Does it describe the risks/benefits of each treatment?" and "Does it provide support for shared decision-making?". The final section (question 16) assesses

the overall rating of the publication ("Based on the answers to all of the above questions, rate the overall quality of the publication as a source of information about treatment choices").

Note: Need to talk about advantages/disadvantages, limitations, relevance to this project etc.

HealthNewsReview

HealthNewsReview provides a set of 10 criteria designed to act as a framework for evaluating the credibility of health-related media. The criteria is based on the various elements that all health-related media should consist of, the criteria is composed of criterion such as "Does the story compare the new approach with existing alternatives?", "Does the story use independent sources and identify conflicts of interest?" and "Does the story appear to rely solely or largely on a news release?".

Note: Need to talk about advantages/disadvantages, limitations, relevance to this project etc.

QIMR

QIMR is a tool developed to monitor the quality of health research reports presented in the media. The tool considers 5 main factors that have been deemed to be correlated with the quality of research reports based on interviews with health journalists and researchers. The 5 main factors are: background information provided, sources of information used, manner in which results were analysed, context of the research and the validity of their methodology.

2.1.2 Studies Applying Checklists to Health Information Online

2.2 Document Classification Methods

The task of document classification is a heavily researched topic due to its wide number of applications in various domains. Formally, document classification is defined to be the task of finding a classifier $f : D \rightarrow L$ where $D = \{d_1, d_2, \dots, d_n\}$ is a collection of documents and $L = \{l_1, l_2, \dots, l_k\}$ is the set of possible labels that a document d can be classified as.

$$f : D \rightarrow L \text{ where } f(d) = l \quad (2.1)$$

Common applications of document classification algorithms are: organization and filtering of news articles, document retrieval, opinion mining, email classification and spam filtering [4].

Traditional machine learning-based approaches for document classification, such as Naive Bayes, Support Vector Machines (SVM) and Random Forests, require the manual extraction of features [4] [6] [11] [13] as they are incapable of performing feature learning

[7]. These features are typically either hand-crafted and domain specific, requiring expert-level knowledge for the task domain [13] or are simple and general but prone to the loss of information (e.g. being unable to account for context) since these features, such as term frequency-inverse document frequency (TF-IDF) scores or bag of words (BoW) models, are unable to represent and account for the sequential nature of text.

2.2.1 Representing Text as Features

Due to the unstructured nature of natural text, it is common practice for natural language processing (NLP) tasks to transform the text into a structured representation that minimises the loss of the information required by the model. However, despite common shallow representation techniques of texts such as BoW, TF-IDF or N-Grams showing respectable performances in the classification of texts [20] [21], these representations are prone to being unable to generalise across multiple domains [14] as the texts are typically represented using simple mechanisms that are unable to represent high-level information such as context or the semantic relationship of words..

TF-IDF Score

The TF-IDF score of a word represents its importance in a specific document relative to all other documents within a corpus. Based on the BoW model, the TF-IDF calculates both the frequency of a word within a document (term frequency) and the

N-Grams

Word2Vec

Word Embeddings

- ngrams
- neural language models e.g. word2vec

2.2.2 Machine Learning Methods for Document Classification

Naive Bayes Classifier

Naive Bayes is a probabilistic classifier that functions on an assumption in how the data (i.e. the words in a document) is generated. The assumption that these Bayesian classifiers are based on is that the distribution of different words within a corpus are independent from each other. Despite this assumption clearly being wrong when considering the distribution of words within a document (due to the sequential nature of text), Naive Bayes classifiers are still able to perform well in document classification tasks such as the filtering of spam emails [15].

Naive Bayes classifiers utilize Bayes' theorem, which attempts to find the probability of an event B occurring given some prior condition A i.e. $P(B|A)$. In the context of

document classification, Naive Bayes classifiers classifies a document d by calculating the probabilities that the document belongs for all labels $l_i \in L$ and then selecting the highest probability [6]:

$$P(L = l_i|d) = \frac{P(l_i)P(d|l_i)}{P(d)} \quad (2.2)$$

According to Aggarwal et al. [4], there are two main types of Naive Bayes classifiers that are commonly used, they are the multivariate Bernoulli model and the multinomial model. The main discriminating factor between these two models is that the Bernoulli model does not take into account the frequency of words as it represents a document using a vector of binary features which signify the presence or absence of words, based on some vocabulary, for a given document. This is in contrast with the multinomial model, which accounts for the frequency of words as the document is represented using a BoW model. Deciding on which model to use for document classification largely depends on the size of the vocabulary as, according to McCallum et al. [12], the multinomial models almost always outperforms the Bernoulli models if the size of the vocabulary is large (> 1000) or even if the size has been optimally chosen for both models.

Support Vector Machines

SVMs are a type of supervised machine learning-based binary classifiers that are extensively used for document classification due to their capability of handling the high dimensional and sparse nature of the common techniques used to represent text documents [10]. SVMs are able to classify a document by using a hyperplane to separate the different classes into separate regions. The hyperplane used is determined by choosing the one that maximises the margin of separation (i.e. The Euclidean distance between the hyperplane and all data points in the representation space) between the two classes [4].

Consider the illustration shown in figure 2.1 that presents three two-dimensional hyperplanes, A , B and C , which separates the classes 'x' and 'o'. Visually, we can determine that A is the hyperplane that maximises the margin of separation, thus, A will be used as a decision rule to classify any new document based on its location with respect to A . Mathematically, determining which hyperplane to use as the decision rule is an optimisation problem that attempts to maximise the margin represented as shown in Equation 2.3:

$$\text{maximise: } \frac{1}{2} \left\| \frac{1}{w} \right\|^2 \quad (2.3)$$

This optimisation problem is often re-framed to the minimisation of Equation 2.4:

$$\text{minimise: } \frac{1}{2} \|w\|^2 \quad (2.4)$$

This technique of determining the hyperplane requires that the two classes are linearly separable. In situations where this isn't true, the kernel trick [5] is applied which is essentially a function that maps data in a particular representation space to another

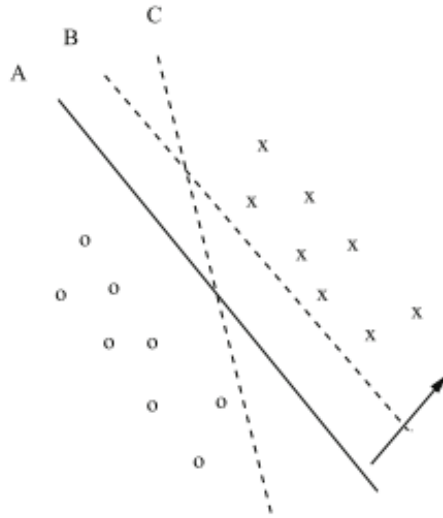


Figure 2.1: What is the Best Separating Hyperplane? [4]

representation space of a different dimension. This change in dimensionality can allow the classes to become linearly separable and thus a hyperplane can then be constructed that separates each class as illustrated in Figure 2.2.

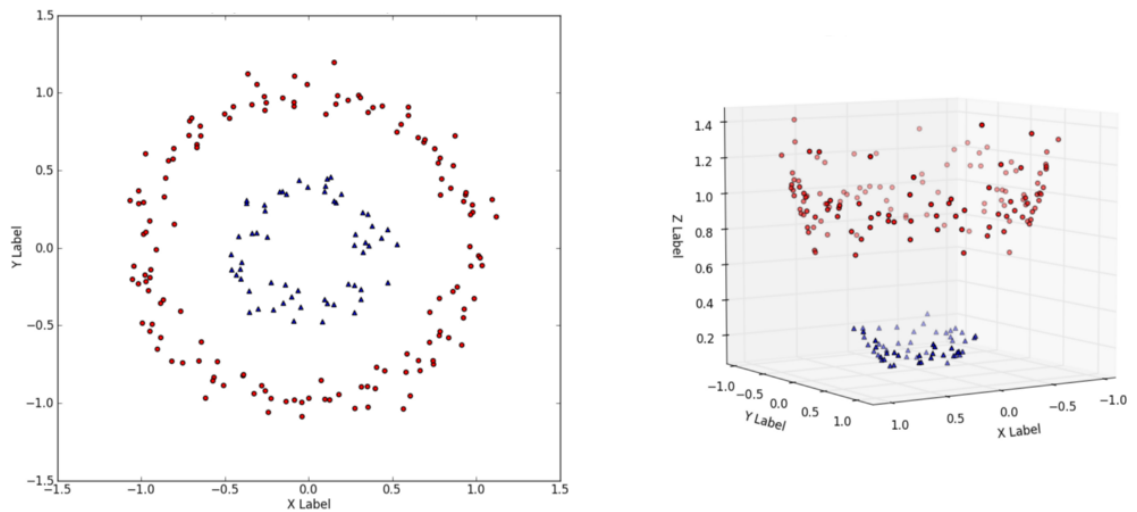


Figure 2.2: (Left) Non-linearly distributed data of two classes in a 2D representation space.
(Right) Linearly separable data of the same classes in a 3D representation space after the application of a kernel [3].

2.2.3 Deep Learning Methods for Document Classification

Deep learning models are a class of machine learning models that have the capability of automatically learning a hierarchical representation of data [7]. These hierarchical representations are constructed through the use of artificial neural networks, the main underlying mechanism of deep learning models. Typically, large amounts of training data is required to train a model in learning the language model required to attain state of the art results, in the task of document classification for instance, the size of commonly used non-domain specific datasets range from hundreds of thousands of training examples to millions [9] [20] (*note: look into the datasets used by state of the art approaches*). Due to these constraints, it is not feasible to procure a dataset for the domain specific task of this project due to the aforementioned knowledge expertise and time requirements to manually label the articles required. Hence, (*Talk about transfer learning/N-shot learning/domain adaptation here*) will be used to overcome this issue.

Introduce the state-of-the-art DL based approaches for document classification and try to compare it performance with state-of-the-art ML approaches

For each model, talk about the following:

NOTE: REMEMBER WHEN WRITING THIS SECTION TO ALWAYS CONSIDER HOW IT DIFFERS TO ML TECHNIQUES

- *How it works and the mechanisms involved*
- *Advantages*
- *Limitations*

Introduce the typical architectures used for document classification e.g. RNNs, LSTMs, CNNs, GRUs?

2.2.4 Deep Learning Models

Recurrent Neural Networks

Gated Recurrent Unit Networks

Long Short-Term Memory Networks

Convolutional Neural Networks

2.2.5 Classifying Documents with Limited Labelled Examples

Transfer Learning

Talk about transfer learning and how it works and how it is applicable to this project.

N-Shot Learning

Talk about zero/few/etc-shot learning and how it works and how it is applicable to this project.

2.3 Discussion

Summarize lit review and describe why DL-based approaches should be preferred over ML-based for this type of problem. Also talk about Transfer/N-Shot learning and describe which one will be feasible given the project's time constraints

2.3.1 Scaling the Assessment of Quality and Credibility of Health Information

Summarise the idea that the previous studies applying instruments/tools to assess the quality of online health information have been relatively small because they rely on experts. It would be useful to be able to develop a way to automate the assessment of online health information because it could lead to new ways of dealing with the spread of misinformation online, and that is hasn't been done properly before.

2.3.2 Using Deep Learning Methods to Automate Quality and Credibility Assessments

Explain that deep learning is a good choice for building an automated assessment tool because it should outperform traditional methods of document classification, but because DL methods usually require large datasets of labelled examples than it makes sense to investigate transfer learning/n-shot learning.

Chapter 3

Proposed Approach

3.1 Rationale

Introduce and discuss the factors that led to me choosing the proposed approach

3.2 Credibility Criteria

Introduce and discuss the 7 criteria that will be classified and describe how the criteria was determined

3.3 Study Data

Talk about the data I'll be using, how we got it, its characteristics etc.

3.4 System Model

Describe the architecture of the model

3.5 Experiments

Describe the experiments that I'm planning to do (in such a way that they are easily reproducible)

3.6 Outcome Measures

Talk about the type of analyses that I'll be doing to determine the performance of my proposed model

Chapter 4

Conclusions and Future Work

4.1 Conclusions

The end

4.2 Future Work

Chapter 5

Abbreviations

AWGN	Additive White Gaussian Noise
BC	Broadcast Channel
BS	Base Station
CSI	Channel State Information
CSIR	Channel State Information at Receiver
CSIT	Channel State Information at Transmitter
dB	Decibels
DPC	Dirty Paper Coding
GS	Gram-Schmidt
RVQ	Random Vector Quantisation
SISO	Single Input Single Output
SNR	Signal to Noise Ratio
SINR	Signal to Interference plus Noise Ratio
MISO	Multiple Input Single Output
SIMO	Single Input Multiple Output
MIMO	Multiple Input Multiple Output
MMSE	Minimum Mean Square Error
MRC	Maximum Ratio Combining
QoS	Quality of Service
TDD	Time Division Duplex
FDD	Frequency Division Duplex
ZF	Zero-Forcing
ZFBF	Zero-Forcing Beamforming
ZMCSCG	Zero Mean Circularly Symmetric Complex Gaussian

Appendix A

name of appendix A

A.1 Overview

here is the Overview of appendix A ...

A.2 Name of this section

here is the content of this section ...

Appendix B

name of appendix B

B.1 Overview

here is the Overview of appendix B ...

B.2 Name of this section

here is the content of this section ...

Bibliography

- [1] “DISCERN - The DISCERN Instrument.” [Online]. Available: <http://www.discern.org.uk/discern{-}instrument.php>
- [2] “Our Review Criteria - HealthNewsReview.org.” [Online]. Available: <https://www.healthnewsreview.org/about-us/review-criteria/>
- [3] “Understanding the kernel trick. – Towards Data Science.” [Online]. Available: <https://towardsdatascience.com/understanding-the-kernel-trick-e0bc6112ef78>
- [4] C. C. Aggarwal and C. Zhai, “A SURVEY OF TEXT CLASSIFICATION ALGORITHMS,” 2012. [Online]. Available: <http://alias-i.com/lingpipe/>
- [5] M. A. Aizerman, E. A. Braverman, and L. Rozonoer, “Theoretical foundations of the potential function method in pattern recognition learning,” *Automation and Remote Control*, vol. 25, pp. 821–837, June 1964.
- [6] M. Allahyari, S. Pouriyeh, M. Assefi, S. Safaei, E. D. Trippe, J. B. Gutierrez, and K. Kochut, “A Brief Survey of Text Mining: Classification, Clustering and Extraction Techniques,” Tech. Rep., 2017. [Online]. Available: <http://en.wikipedia.org/wiki/Statistics>
- [7] I. A. Basheer and M. Hajmeer, “Artificial neural networks: fundamentals, computing, design, and application,” Tech. Rep., 2000. [Online]. Available: www.elsevier.com/locate/jmicmeth
- [8] D. C. Burgess, M. A. Burgess, and J. Leask, “The MMR vaccination and autism controversy in United Kingdom 1998-2005: Inevitable community outrage or a failure of risk communication?” *Vaccine*, vol. 24, pp. 3921–3928, 2006. [Online]. Available: <https://ac.els-cdn.com/S0264410X06002076/1-s2.0-S0264410X06002076-main.pdf?{-}tid=46d1dda6-f576-4f5e-ad53-d550f1cd9990{&}acdnat=1534726962{-}1b237371d8bb916694f34f0f951c84bc>
- [9] A. Conneau, H. Schwenk, Y. Le Cun, and L. Loic Barrault, “Very Deep Convolutional Networks for Text Classification,” 2017. [Online]. Available: <https://arxiv.org/pdf/1606.01781.pdf>

- [10] F. Informatik, “UNIVERSIT AT D ORTMUND Text Categorization with Support Vector Machines: Learning with Many Relevant F eatures Thorsten Joachims,” Tech. Rep., 1997. [Online]. Available: https://www.cs.cornell.edu/people/tj/publications/joachims_{_}97b.pdf
- [11] V. Korde and C. N. Mahender, “TEXT CLASSIFICATION AND CLASSIFIERS: A SURVEY,” *International Journal of Artificial Intelligence & Applications (IJAIA)*, vol. 3, no. 2, 2012. [Online]. Available: <http://aircconline.com/ijaia/V3N2/3212ijaia08.pdf>
- [12] A. McCallum and K. Nigam, “A Comparison of Event Models for Naive Bayes Text Classification,” Tech. Rep., 1998. [Online]. Available: <http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.65.9324{&}rep=rep1{&}type=pdf>
- [13] K. Pasupa, “A Comparison between Shallow and Deep Architecture Classifiers on Small Dataset,” 2016.
- [14] R. Rosenfeld, “TWO DECADES OF STATISTICAL LANGUAGE MODELING: WHERE DO WE GO FROM HERE?” Tech. Rep., 2000. [Online]. Available: <https://www.cs.cmu.edu/{~}roni/papers/survey-slm-IEEE-PROC-0004.pdf>
- [15] M. Sahami, S. Dumais, D. Heckerman, and E. Horvitz, “A Bayesian approach to filtering junk e-mail,” Tech. Rep. Cohen, 1998. [Online]. Available: <http://research.microsoft.com/en-us/um/people/horvitz/junkfilter.htm>
- [16] S. Sommariva, C. Vamos, L. U.-L. Mantzarlis, Alexios Dao, and D. Martinez Tyson, “Spreading the (Fake) News: Exploring Health Messages on Social Media and the Implications for Health Professionals Using a Case Study,” *American Journal of Health Education*, vol. 49, no. 4, pp. 246–255, jul 2018. [Online]. Available: <https://www.tandfonline.com/doi/full/10.1080/19325037.2018.1473178>
- [17] “Germany investigating unprecedented spread of fake news online — World news — The Guardian,” 2017. [Online]. Available: <https://www.theguardian.com/world/2017/jan/09/germany-investigating-spread-fake-news-online-russia-election>
- [18] S. Vosoughi, D. Roy, and S. Aral, “SI:The spread of true and false news online,” Tech. Rep. 6380, 2018. [Online]. Available: <http://science.sciencemag.org/content/sci/359/6380/1146.full.pdf{ }0Ahttp://www.sciencemag.org/lookup/doi/10.1126/science.aap9559>
- [19] D. Zeraatkar, M. Obeda, J. S. Ginsberg, and J. Hirsh, “The development and validation of an instrument to measure the quality of health research reports in the lay media,” *BMC Public Health*, vol. 17, no. 1, p. 343, dec 2017. [Online]. Available: <http://bmcpublichealth.biomedcentral.com/articles/10.1186/s12889-017-4259-y>
- [20] X. Zhang, J. Zhao, and Y. Lecun, “Character-level Convolutional Networks for Text Classification,” 2015. [Online]. Available: <https://arxiv.org/pdf/1509.01626.pdf>

-
- [21] C. Zhou, C. Sun, Z. Liu, and F. C. M. Lau, “A C-LSTM Neural Network for Text Classification,” 2015. [Online]. Available: <https://arxiv.org/pdf/1511.08630.pdf>