

# Topical Rumor Detection based on Social Network Topic Models Relationship

Diogo Nolasco<sup>1</sup>, Jonice Oliveira<sup>1</sup>

<sup>1</sup>Programa de Pós-Graduação em Informática – Universidade Federal do Rio de Janeiro (UFRJ)  
Rio de Janeiro, RJ – Brazil

diogo.sousa@ppgi.ufrj.br, jonice@dcc.ufrj.br

**Abstract.** *The rumor detection problem on social networks has attracted considerable attention in recent years with the rise of concerns about fake news and disinformation. Most works focused on detecting rumors individually, classifying whether a post is a rumor or not. This paper proposes a topic-level method for rumor detection. We use a topic model method on social and scientific domains and correlate the topics found to detect the most prone to be rumors. Two scenarios were analyzed; the Zika epidemic and the Brazilian presidential speeches. Results of the former suggests that the least correlated topics contains rumors and local discussions. The latter suggests a strong correlation of rumor topics from both speeches and social domains.*

**Keywords.** *Text Mining; Topic Modeling; Social Networks; Topic Labeling; Topic Correlation*

## 1. Introduction

The traffic and discussions generated on social networks have been increasing with their development over time. In recent years, breaking news appears first on microblogs, before making it through to traditional media. A discovery, event, or any information can become viral almost instantly. While the quantity generates massive data for analysis, the quality of information does not become better. All kinds of false information, especially rumor information, have acquired an unprecedented range and permeates most social communities. Consequently, the means of automatically detecting the information's credibility and monitor public subjects has been getting increased attention.

Rumor detection is one of the research topics critical to social networks. Rumor itself is often viewed as a tale of explanations of events circulating from person to person and pertaining to an object, event, or issue in public concern [Peterson and Gist 1951]. With the massive amount of data in social networks, it is hard to distinguish reliable information from false information. The rumor diffusion can happen inadvertently or maliciously, so, commonly, the message appears to be truthful. The difficulty of the task is so that many news agencies and organizations have departments to assess and inform the public about false information vehiculated on social media. This spread can cause people to make wrong or misinformed decisions and could even harm social stability. Thus, its detection is a major concern for social networks and society.

There are a number of studies on rumor detection, with most of them basically consisting of the task of classifying a message into reliable or unreliable (binary classification). However, there are few studies that try to analyze rumor subjects and detect rumor topics at a coarse-grained level [Cao et al. 2018].

In this paper, we study the problem of automatically detect rumor topics spreading in social media. We propose a method for detecting rumors using topic models to find unreliable topics given a main event or subject. Our methods use two datasets, the social network dataset, which possibly contains rumors, and a ground truth dataset, that shares the same main event or subject of the social network one but contains curated, technical reliable data or rumor verified data. With these topics, we are going to use a topic correlation approach to establish relationships between the two datasets at a topic level. Finally, with the topic correlation, we draw some features that suggest which topics contain rumors and discuss the results to bring more evidence. The main contributions of this work are:

- The cross-topic model method for inferencing rumors.
- A topic correlation approach to detect rumor topics.
- An assessment of the nature of rumor topics included in the relationships between different topics.

The rest of the paper is organized as follows; Section 2 covers related works. Section 3 provides background knowledge on rumors and topic models. Section 4 presents the proposal and methods used. Section 5 presents the results followed by an in-depth discussion. Finally, Section 6 concludes this work.

## 2. Related Works

There are limited works that study the rumor detection at a topic level. Several survey papers exist on the subject but none of them presents relevant works on the matter [Cao et al. 2018; Zubiaga et al. 2018]. Some of them address the detection in the context of fake news [Ahsan and Kumari 2019; Sharma et al. 2019].

Although not at the topic level, the traditional methods for rumor detection are related to this work. Basically, related studies focus on extracting useful and efficient features for rumor detection. Generally speaking, features for rumor detection can be divided into three types: (1) content-based features; (2) user-based features; (3) propagation-based features.

For content-based features, [Ratkiewicz et al. 2011] identifies misleading political memes on Twitter using content-based features, including hashtags, links, and mention. [Takahashi et al. 2015] computes the ratio of the number of rumor and non-rumor messages vocabulary words as a feature to detect rumors. They found that vocabulary distributions are different between them.

For the user-based features, [Castillo et al. 2011] used a number of user features like age, number of posts, followers, number of friends, and others to detect rumors. Other works like [Al-Khalifa and Al-Eidan 2011] also used several other user attributes.

For the propagation-based features, [Mendoza et al. 2010] analyzed the retweet network topology and the diffusion patterns of rumors and discovered that they are

different from traditional news, they also found that rumors tend to be questioned more than news by the Twitter community. There is also the work of [Kwon et al. 2013], that discovered that rumor tweets had more cycle volatility, compared with non-rumor tweets.

### 3. Background

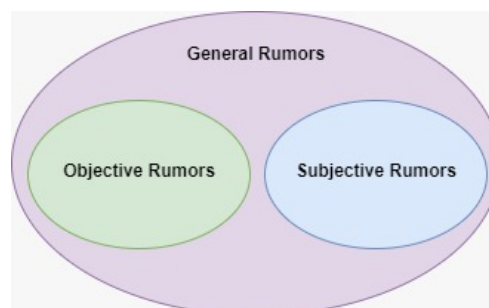
The proposal presented in this work is based on topic modeling techniques that are used to do rumor detection at a topic level, i. e., detecting rumor topics. Thus, an overview of rumors and topic models are presented in this Section.

#### 3.1. Rumor

There are various definitions of rumor in different areas. There is the view that a rumor is a story or statement in general circulation without confirmation or certainty to facts [DiFonzo and Bordia 2007]. While [Allport and Postman 1947], define it as a story or a statement whose truth value is unverified or deliberately false.

The existence of different definitions makes it hard to compare the effectiveness of different methods for rumor detection. However, there are some typical definitions usually found in the literature (as in Figure 1) of this research area[Cao et al. 2018]:

- **General Rumors:** Concerns pieces of information where the truth value is unverified, both inadvertently or on purpose. Gossips, fake news, and unverified information are examples of this broad concept.
- **Objective Rumors:** Rumors that are verified fake information, i.e., where there exist reliable sources that show that the information is false. Messages that are spread or viral by people that do not know the sources or by people that deliberately spread the rumor.
- **Subjective Rumors:** Rumors where the truth value is determined by the subjective judgment of users. For example, where veracity is based on people's subjective feelings.



**Figure 1. Rumor types**

In this work, we study mainly the objective rumors, those that are proven to be false information. These include “fake news”, disinformation and misinformation. This kind of rumor is more easily comparable and can be assessed by using a ground source of truthful values for the information.

### 3.2. Topic Modeling

Topic models are a suite of algorithms used for discovering the abstract "topics" that occur in a large collection of documents through statistical models. Topic modeling is a frequently used text-mining tool for the discovery of hidden semantic structures in a text body [Blei et al. 2010]. It is considered an unsupervised learning technique that learns classes or topics from previously untrained data, effectively finding patterns on text and relations between terms.

As an unsupervised task, there is no need for training data and the classes called topics on this context to emerge from the application of the statistical model on the collection. In the context of social networks, it is usually unknown to the user which topics exist in the data beforehand. The social discussions are fluid and dynamic, so the topics of discussion change over time. The very task of classifying those in categories needs an increased effort of domain specialists. New topics emerge every day in social media and the supervised learning techniques are often not feasible to be applied in this scenario.

In this work, we use topic models to find the topics present in a social network and scientific datasets respectively. It is suitable for the use in social networks as it is difficult to know the discussion topics beforehand, and in the scientific dataset, it provides an overview of the research topics that are researched in the academic community.

For this task, we use the Latent Dirichlet Allocation (LDA) model [Blei et al. 2003]. It is one of the most popular and is the source of many recently created models. The LDA was based on two other popular topic models: Latent Semantic Allocation (LSA) [Steyvers and Griffiths 2007] and Probabilistic Latent Semantic Indexing (pLSI) [Blei and Lafferty 2009].

It works by creating two types of multinomial probabilistic distributions. The first is a distribution of a term over the topics where each term has a probability associated with each topic (which represents the relevance of the term for each topic). The second one is a distribution of a document over the topics (which represents the relevance that each topic has in the document). This second distribution can be found by using the first one as a document is viewed as a collection of terms. Based on these distributions one can group the documents according to the probabilities associated with each topic.

Thus, given a vocabulary  $V = \{w_1, w_2, \dots, w_{|V|}\}$  consisting of all terms that exist in the collection  $C = \{d_1, d_2, d_3 \dots\}$ , where  $w_n$  is the  $n$ th-term of the vocabulary and  $d_n$  is the  $n$ th-document of the dataset, we define a topic model from the algorithm as:

*DEFINITION 1.* A topic model  $\theta$  in  $C$  is a probability distribution of terms such that:

$$\theta = \{p(w_1|\theta), p(w_2|\theta), \dots, p(w_{|V|}|\theta)\} \quad (1)$$

and

$$\sum_{w \in V} p(w|\theta) = 1 \quad (2)$$

A "Zika" topic, for example, would assign higher probabilities to words like "epidemic", "vaccine" and "victims" and lower probabilities to less relevant words such as "fun" and "music". Basically, a topic is a probabilistic distribution where words with

high probability are more relevant to the assigned topic, while low probability words are irrelevant or stop words.

#### 4. Proposal

The objective of the proposal is to find objective rumors from topic correlation obtained from social network and authoritative dataset topics respectively. In this work, we use a scientific dataset and a Brazilian presidential speech dataset as our authoritative dataset, i.e., a dataset whose topics are reliable and verified. These datasets are used respectively in two scenarios. The first is conducted mainly in the Zika epidemic scenario to find rumors that were spread through social media and thus, a scientific dataset provides the topics and labels that could be used to verify information credibility of what is discussed about the disease in the social networks. The second one is executed in a different domain, the Brazilian presidential speeches occurred in 2020 in the context of the COVID-19. Here we analyze how rumors presented in a public speech by an authority relate to social discussions.

The proposal technically consists of three different tasks performed in those datasets: (1) Topic Detection; (2) Topic Labeling, and (3) Topic Correlation. Topic detection consists in extracting topics from collections. Topic Labeling is used to analyze, correlate, and provide a comprehensive overview of the topics extracted. Topic Correlation is the task where the relationships between the two datasets are established and where we discuss the findings concerning rumor detection.

The techniques used for Topic Detection and Labeling are an implementation of the ones described by [Nolasco and Oliveira 2018], as they prove successful and efficient when applied in the dataset types used in our experiments, namely, in the academic and social network domains. The Topic Correlation approach is based on the Kullback-Leibler divergence [Kullback and Leibler 1951], a measure used to associate two probabilistic distributions and that we use on the distributions resulted from topics models.

As the results and findings depend on the methods used, we give a brief explanation of these three tasks in the next sections to provide the necessary technical background to understand our study.

##### 4.1. Topic Detection

In this task, we are using the Latent Dirichlet Allocation (LDA) to detect topics from the textual collections. For this, we also use the topic detection methods described in [Nolasco and Oliveira 2018].

Usually, a topic model algorithm (including LDA) expects two main input parameters, a number of topics to find  $K$  and a collection of  $N$  documents. In the social network domain, the topics are not known by the user and thus the  $K$  value can vary according to the dataset. The topic detection method used here circumvents this by automatically estimating the best  $K$  value for a particular dataset based on Stability Analysis.

After execution, this algorithm results in a series of topics  $\Theta$  where each topic is a set of terms with associated probabilities.

## 4.2. Topic Labeling

While the topics are just probabilistic distributions of terms over the topics, it is essential to the rumor analysis that these distributions could translate into textual themes or subjects for comparison and assessment. For this task, we use the topic labeling method presented in [Nolasco and Oliveira 2016]. It was made to work with topic models with successful results when labeling social and science topics. It consists of three steps illustrated by Figure 2, they are: 1) Candidate Label Selection, where keywords and keyphrases are extracted from the topic relevant documents; 2) Score and Ranking, where the keywords are ranked based on a function score and 3) Label Selection. Where the labels are attributed to each topic.

## 4.3. Topic Correlation

For comparing the potentially unreliable social topics and the scientific truth source topics, we are using the correlation between the two topics as a means to validate the topic's reliability. We calculate it using a distance function based on the KL-Divergence.

The KL-Divergence is one of the most used when comparing two probability distributions (in this case the topics) but it is not asymmetric measure, i.e., inverting the parameter order results in different metrics values. Since the topic correlation needs to be the same on both directions we use a distance based on that divergence defined by its creators [Kullback and Leibler 1951] and calculated as:

$$\text{KL-Distance}(\theta_1, \theta_2) = D_{KL}(\theta_2 || \theta_1) + D_{KL}(\theta_1 || \theta_2) \quad (3)$$

Which is symmetric and nonnegative and where  $D_{KL}$  is the original KL-Divergence between two topics. A low total value means that little information is lost when comparing the two distributions while a high value means that they are very dissimilar.

For comparison between topics and terms that are not represented by a probability distribution, i.e., the public speeches used in this work, we use another approach based on the overlap similarity measure [Vijaymeena and Kavitha 2016], as follows:

$$\text{overlap}(X, Y) = \frac{|X \cap Y|}{\min(|X|, |Y|)} \quad (4)$$

Where X and Y are sets, in our case X and Y are topics represented here as a set of terms. We chose this similarity measure because we are comparing different sized finite sets and the public speech (represented by X) set was manually built by an expert and thus contains a more curated list of terms. Meanwhile, social conversations sets (represented by Y) were automatically generated by the topic detection algorithm and may contain some irrelevant words or noise. To highlight the importance of the president's speech terms, we weighted the intersection.

To calculate this similarity we use all terms available from the president's speeches sets and the 100 most relevant terms from our social topics, since the speech topics are relatively small and the social topics contain a vocabulary of thousands of terms. A preprocessing step consisting of removing links, emojis and stemming words is also applied to data to favour matches of derivative words such as “argues” and “argued”. This metric ranges from zero to one where a zero means that the sets are

distant or unrelated while one means that  $X$  is a subset of  $Y$ . This way we can effectively measure the anticipation and impact between official speeches and social discussions.

The difference between reference or authoritative sets and between correlation methods gives a measure of how we can compare sets of topics from different perspectives and find rumors by correlation, even when using different domains or metrics.

## 5. Results

We conduct two experiments to evaluate how rumor topics can be found by comparing social and reference topics.

The first uses the Zika epidemic scenario to showcase how the topic correlation between social network and science topics can be used to detect unrelated or unreliable information about the event and the disease itself. We also use the topics to differentiate general rumors, local discussions, and fake news.

The second compares Brazilian presidential speeches in the context of COVID-19 by president Jair Bolsonaro in 2020 to show how we can use a set of terms as topics from official sources to detect rumors and analyze their influence on society.

The first evaluation (Zika epidemic) is made using two datasets, a Twitter dataset of Zika related posts and a PubMed corpus of Zika related articles. The scenario covered by these datasets is relative to the context of the Zika epidemic from 2015 to 2016, which contains a variety of topics like reports, propagation to various countries, associated diseases, and influence on the 2016 Olympic Games organization.

The second evaluation (Presidential speeches) is made with the aid of specialists who aided in the identification of rumors and disinformation on Jair Bolsonaro's public speeches occurred on the first half of 2020. The speeches were separated in a set of related terms as "manual topics" where the specialists pointed what topics from the discourse were related with disinformation.

A qualitative study is conducted to analyze the social topics and its correlation to the science topics and speech topics as a means to detect rumors about the disease (in the Zika epidemic scenario) and rumors about social distancing and the COVID-19 (in the Presidential speeches scenario). Since we work with objective rumors, we are using some "ground truth" data to support our claims and findings. Specifically, we are going to use the timeline of the Zika epidemic communications report [Fundação Oswaldo Cruz. and Araujo 2016], the news reported by media for social topics, the manual analysis of specialists (in the Presidential speeches scenario), and the World Health Organization list of rumors about the Zika virus [WHO 2016] and the Covid-19 [WHO 2020]. The timeline contains the main events reported to the public such as the virus identification and the outbreak declaration by authorities, while news can be used to verify associated secondary events and subjects such as local efforts to combat the epidemic and ways of preventing the infection. The list of rumors contains a compilation of the all the false information spread or divulged during the outbreak and can be used to verify reliability.

## 5.1. Datasets

Two databases were used for the experiments corresponding to the two different scenarios, i.e., the Zika epidemic and the Brazilian presidential speeches.

The database for the first experiment was made extracting posts from around the world with the #zika “hashtag” and articles from the PubMed database containing the keyword Zika. The term is popular in both domains and has little ambiguity, so the addition of other terms could introduce more noise to the data. The time span of the datasets covers documents created from May. 2015 to Dec. 2016, to cover the lifespan of the epidemic. A total of 85,601 tweets and 1769 articles were retrieved. A preprocessing was made in these data by removing emotes, links, and accents from the text.

For the second experiment, to analyze the public perceptions and draw parallels with the Brazilian presidential speeches, we used a dataset of Twitter posts in Portuguese associated with Covid-19 and the transcriptions of president Bolsonaro’s official speeches. The transcriptions were taken from the speech dates, respectively: March 12th, March 24th, March 31th, April 8th and April 16th. Tweets were acquired in the days after each speech, being collected on the following dates: March 13th, March 25th, April 1st, April 9th and April 17th. We chose these dates to get data from social discussions occurring after each speech as a means to track rumors and repercussions of the pronouncement. We opted for this time frame in order to minimize noise in our data, since various other events could influence social discussions. The dataset comprises a total of 15,154 tweets, which includes posts and replies, and excludes retweets. We searched for tweets using the terms: "coronavirus", "corona virus", "corona viruses", "coronaviruses", "viruscorona", "virus corona", "sars virus", "sars-like virus", "virus sars", "sars viruses", "sars-like viruses", "Wuhan virus", "HCoV-229E", "HCoV-OC43", "SARS-CoV", "HCoV-NL63", "HCoV-HKU1", "MERS-CoV", "HCoV-EMC", "2019-nCoV", "2019nCoV", "ncov19", "ncov-19", "covid-19", "covid19", "sars-cov-2", "Severe acute respiratory syndrome", "Middle East respiratory syndrome", "Wuhan pneumonia". We also used the option “lang: pt” to collect messages preferably in Portuguese. Just like the previous database, a preprocessing was made by removing emotes, links, and accents from the text

## 5.2. Execution

Since the databases and results are in different contexts and may have different conclusions, we divide the execution in two parts concerning the two scenarios. First, the Zika epidemic shows the results in a public health context and secondly, the Brazilian presidential speeches showing a more political context amidst the rumors involving the Covid-19 pandemic such as social distancing and the use of medicines without scientific proof.

### 5.2.1. Zika Epidemic

The proposal consists of three tasks, Topic Detection, Labeling, and Correlation. In this Section, we present the details of each task execution in the scenario of the Zika epidemic.

For the topic detection and topic labeling, we used the same parameters that [Nolasco and Oliveira 2020] used on similar datasets that provided good results.



Specifically, the K number of topics varying between 4 and 20 and the top 10 documents and words for D and W parameters in their candidate selection algorithm. For the topic correlation, we use the symmetric KL-Divergence that is a measure of similarity between the topics from the two datasets and does not need any special parameter.

In this experiment, we analyzed two different periods of research and social discussion: 1) From May. 2015 to Feb. 2016, covering the start of the epidemic and first counter-measures and 2) From Mar. 2016 to Dec. 2016, covering the apex of the outbreak and the subsequent decline.

Table 1 shows results extracted from social media in the two periods while Table 2 shows the research topics extracted during the same periods. We show only the relevant topics found in each dataset, i. e., the topics associated with themes present in news sources, the Zika timeline, or the WHO documents. This resulted in 8 social topics and 5 science topics for the first period; 8 social topics and 7 science topics for the second period.

**Table 1. Social Topics for Zika epidemic**

Period 1 – From May. 2015 to Feb. 2016	
Topic	Labels
1	zika virus’ — doctors expose monsanto linked pesticide, birth defect microcephaly, birth defect
2	zika virus #zikainrio #zika virus @rio2016 en, cancelling rio olympics due, skipping #2016olympics due
3	world health organization director general declares #zika virus outbreak, world health organization declares spread, intl health regulations emergency committee
4	miami beach #zika virus #zikazone #advisory #miamibeach, caution pregnant women advised, #miami #beach area
5	prevenir el #zika #zika virus pandemia ubicada como peligro mundial hoy @hijoslakebuena, si estas embarazada redobla el cuidado contra el mosquito del dengue, #zika virus el virus zika es causado por la picadura de
6	zika vaccine candidates #zika #zika virus #cdc #nih #niaid #vaccines \$gsk \$sny, zika vaccine candidates #zika #zika virus #cdc #nih #niaid #vaccines \$sny \$gsk, zika \$nlk #zika #zika virus #vaccines #pharma #nih #cdc \$sny \$gsk
7	zika virus spreads #zika virus #automotive #india, zika virus spreads, risk low
8	caso de, primeiro caso, zika vírus
Period 2 – From Mar. 2016 to Dec. 2016	
Topic	Labels
1	neutralizing human antibodies prevent #zika virus #zikv replication, human protein ifitm3 blocks #zika virus replication, human fetal neural stem cells
2	2016, transmission, #cdc, sexual, cdcgov
3	fight #zika #doyourjob @housegop @senategop #zika virus, fight #zika virus ravaging fl, fighting #zika virus fails
4	#nc governor pat mcrrory, dilemma, #miamibeach
5	#cuba reports 1st #zika travel case, #breaking beijing reports 3rd case, chp confirms #zika virus case
6	asian zika virus mutated negatively & zika virus mutated negatively & zika virus mutated negatively
7	mosquito repellent zika virus protection, 99 free ship
8	#zika virus, cientistas #vooz, #vooz #zika virus, solucoes baseadas em #dados para fazer frente ao #zika virus

**Table 2. Science Topics for Zika epidemic**

Period 1 – From May. 2015 to Feb. 2016		Period 2 – From Mar. 2016 to Dec. 2016
Topic	Labels	Labels
1	ZIKV, virus, infection	Zika virus prevention, travellers concern, emerging infectious diseases

2	Zika, emerging doorstep, outbreak	Congenital fetal malformations, pregnant women, congenital microcephaly
3	Brazil, Bahia, Americas	Zika virus infection, emergency department, ZIKV IgM
4	Following dengue, dengue spread, zika	Dengue Virus, human semen, pregnant
5	Zika virus infection, co-infection, new threat	Counter zika virus, diagnostic challenge, detecting
6		Neurologic inhibition, inflammatory, imported arbovirus
7		Mosquito-borne arboviruses, African, saliva

Next, we use a topic correlation on these topics to calculate a similarity between them to find the reliability of the information. Specifically, we use the symmetric KL-Distance to measure the distance between social and science topic distributions. For this analysis, we calculate the correlation between social and science topics from the same period.

Table 3 shows the results of the two periods, both as a heat map. Results are truncated with a precision of two. The cell colors of the heat map are such that the green indicates the lowest distances or most similar while the red indicates the most distant in comparison with the others and to facilitate interpretation.

**Table 3. KL-Distance between social and science topics over two periods**

Period 1 – From May. 2015 to Feb. 2016										
Social Topics	Science Topics							Avg		
		1	2	3	4	5				
	1	4.02	2.28	4.71	1.03	2.48	2,90			
	2	2.44	1.99	1.33	1.47	1.98	1,84			
	3	1.21	1.34	3.17	4.93	1.61	2,45			
	4	2.86	3.96	4.41	2.87	1.09	3,04			
	5	4.58	3.98	2.71	4.13	2.75	3,63			
	6	1.46	4.30	3.14	1.47	2.16	2,51			
	7	0.74	2.34	0.67	2.71	1.72	1,64			
8	3.44	4.26	4.92	3.74	2.80	3,83				
Period 2 – From Mar. 2016 to Dec. 2016										
Social Topics	Science Topics							Avg		
		1	2	3	4	5	6		7	
	1	4.90	4.46	2.68	0.52	2.11	4.71		1.38	2,97
	2	4.61	1.12	3.55	0.70	3.75	4.22		2.14	2,87
	3	2.44	4.60	4.87	0.92	3.01	0.54		2.49	2,70
	4	2.60	2.18	2.65	2.50	2.97	2.89		2.61	2,63
	5	0.51	3.43	1.37	3.98	4.49	1.30		2.65	2,53
	6	3.13	4.50	3.96	1.93	3.92	1.41		3.27	3,16
	7	1.18	1.11	1.11	4.25	3.03	3.43		3.59	2,53
8	3.54	3.59	4.70	2.45	2.84	3.84	4.01	3,57		

### 5.2.2. Brazilian Political Speeches

The same set of tasks, Topic Detection, Labeling, and Correlation, is used in this scenario. The differences here is the presence of a set of curated terms representing the subjects of each speech and the use of the overlap metric in the correlation task. Here we describe these with further detail.

To start, in this experiment we are comparing speech topics with social topics. For the speech topics we had an aid of two specialists in the domain to separate the speech in subjects and select representative terms to act as manually annotated topics. They classified each subject in regard of the presence of disinformation, which we treat like the main rumors in this scenario. The discussion and the rumors are based in the WHO advices against rumors too [WHO 2020]. This was done because of the amount of speech data available, which comprises a limited vocabulary and few words, making it not suitable for topic modeling execution. For the social topic detection and labeling, we used the same parameters and process used in the previous experiment. Specifically, the K number of topics varying between 4 and 20 and the top 10 documents and words for D and W parameters in their candidate selection algorithm.

The topic correlation was done using the overlap metric because the speech topics are not a probability distribution, thus the KL-Divergence is not suitable to use in this case. This metric was chosen because we have to compare a curated list of few words related to a topic with a topic model applied on a large set of data. The overlap reinforces the term intersection, which is ideal for this case when we know that the speech terms were chosen by humans.

In this experiment, we analyzed five different speeches and the subsequent social discussions that emerged from it. The speeches comprise the following days respectively: March 12th, March 24th, March 31th, April 8th and April 16th. The tweets were taken accordingly from the subsequent days of each speech: March 13th, March 25th, April 1st, April 9th and April 17th.

Table 4 shows the topics manually annotated from the speech data. It contains a human label to facilitate comprehension and the disinformation field according to the presence or not of rumors in the discourse given by the specialists. On March 12<sup>th</sup> [BrasilGov 2020a], Jair Bolsonaro gave his first presidential official speech addressing the coronavirus health emergency. On this occasion, the president recognized the pandemic crisis and the health risks involved, especially for the elderly. The president did not disseminate any misleading claims on this occasion according to the specialists.

On March 24th [BrasilGov 2020b], despite defending that the virus would quickly go away and that there was no need to follow social distancing steps, such as working from home, Bolsonaro praised Brazil's then Minister of Health, Luiz Henrique Mandetta. He criticized the press and accused the media of spreading fear and hysteria throughout the country. He also referred to "a medicine for malaria", namely chloroquine (also called hydroxychloroquine), as a promising cure for the disease. Topics associated with disinformation at this time include the accusations of media outlets and politicians spreading panic (Topics B and E), hydroxychloroquine use (Topics H and I) and people affected by the virus (Topics F and G).

On March 31th [BrasilGov 2020c], the president argued that social distancing interventions should not be applied to society as a whole and praised

hydroxychloroquine, although he acknowledged that there was no proof of its therapeutic effectiveness. Economic initiatives were announced, such as financial assistance and expanded business credit. He praised the army and the importance of keeping the economy going and safeguarding jobs. To promote normality, to justify the trade-off between health and the economy (Topic B), and to laud the virtues of hydroxychloroquine against coronavirus, he took advantage of false claims (Topic F).

On April 8th [BrasilGov 2020d], he blamed ministers who "disagreed with him" and accused governors and mayors of enforcing needless orders to remain at home. He reinforced the economic cost of the steps of social distancing, insisted that they could not last much longer and articulated the value of the economy as a counterpoint to the epidemic. The president once again disseminated disinformation about hydroxychloroquine (Topic D), insisting it would be remembered as a medicine that saved thousands of lives in Brazil.

Finally, on April 16th [BrasilGov 2020e], he gave a speech after firing Luiz Henrique Mandetta, then Minister of Health, announcing Nelson Teich as the new minister. For causing hysteria and exaggerating the health crisis in the country, he criticized governors and the media. Bolsonaro stressed the dichotomy between public health and the economy of the country, once again addressing the need for quarantine to be abandoned. When arguing that social distancing steps should be relaxed, he again used misleading arguments (Topic D).

**Table 4. Speech Topics by Official Public Speech**

Topic	Human Label	Disinformation
<b>Public speech 1 (March 12nd)</b>		
A	Disincentive to political protests for recognition of the pandemic	No
<b>Public speech 2 (March 24th)</b>		
A	Preparation/ Strategy of the public health system	No
B	Importance of Panic Containment	No
C	Media spread fear / Brazil will be less affected by the virus	Yes
D	Incentive to economy / Normality	No
E	Speech against state and municipal authorities	Yes
F	The only groups at risk are the elderly	Yes
G	Healthy people will not be affected by the virus	Yes
H	Chloroquine effectiveness	Yes
I	Researchers are trained to cure the disease	Yes
J	Tribute to health professionals	No
<b>Public speech 3 (March 31st)</b>		
A	Concerns about health and economy	No

B	Exaltation of maintaining work / economy	Yes
C	Informal professionals / Access to health care	No
D	Economic aid programs	No
E	Elderly / Unemployment	No
F	Chloroquine's proven effectiveness	Yes
G	Minimizing human losses in the face of economic losses	No
H	G20 decisions	No
I	Operation of the armed forces in the pandemic context	No
J	Chloroquine production in the country	No
K	Economic recovery	No
L	Tribute to essential service professionals	No
<b>Public speech 4 (April 8th)</b>		
A	Criticizes the performance of ministers	No
B	Governors and mayors are solely responsible for stay-at-home orders	Yes
C	Exaltation of the importance of the economy in counterpoint to the disease	Yes
D	Medical oath of the effectiveness of chloroquine for curing COVID-19	Yes
E	Economic program with emergency assistance	No
F	Importance of resuming work and economy	No
<b>Public speech 5 (April 16th)</b>		
A	Resignation of the Minister of Health (Mandetta)	No
B	The media create a climate of hysteria hindering health and the economy	No
C	Discussion about the future minister (Nelson Teich)	No
D	Social distancing violates constitutional law	Yes
E	Process of changing ministers	No

Table 5 shows the social topics extracted from social media after each speech, along with the label for comprehension and the top terms of the topic. In total we have 1 speech topic and 9 social topics relative to the March 12<sup>th</sup> speech; 10 speech topics and 4 social topics relative to March 24<sup>th</sup>; 12 speech topics and 4 social topics relative to March 31<sup>th</sup>; 6 speech topics and 7 social topics relative to April 8<sup>th</sup> and 5 speech topics and 5 social topics relative to April 16<sup>th</sup>.

**Table 5. Social Topics for Brazilian Political Speeches by Public Speech**

Topic	Human Label
-------	-------------

<b>Post-public speech 1 (March 13)</b>	
1	Virus severity and people's fears
2	China's relation with the development of the virus, possible medicines and vaccines
3	Political manifestations scheduled for March 15 and confirmed cases worldwide
4	Local protective measures
5	Memes and jokes about the virus
6	The importance of quarantine and social distancing
7	Confirmed cases among close friends and family
8	Suspension of sports events, work and school
9	Rumor that the president was infected due to a declaration of his son
10	Encouragement on prevention measures
<b>Post-public speech 2 (March 25)</b>	
1	Discussions minimizing the severity of the pandemic and criticism of political figures that were supposedly taking advantage of the crisis
2	Jokes about the infection of Prince Charles and an excerpt of the presidential speech in which Bolsonaro claimed military and athletes have minor symptoms when infected
3	Reference to Bolsonaro's comparison of the Coronavirus with a "little flu"
4	Defense of the president's positions
5	Authorization of hydroxychloroquine prescription for Covid-19 cases
<b>Post-public speech 3 (April 01)</b>	
1	Criticism of Bolsonaro's political opponents
2	Criticism of media outlets
3	The president's decisions during the pandemic (hydroxychloroquine import, military help and financial aid measures)
4	General information about the virus and related health behavior
<b>Post-public speech 4 (April 09)</b>	
1	Discussions about quarantine, its feasibility and effect among the poorest
2	Criticism of containment measures and complaints about pandemic
3	Criticism of media coverage and legacy media journalists
4	Social distancing fatigue
5	Defense of the use of hydroxychloroquine as a medicine for Covid-19
6	Criticism of the minister of health and WHO

7	Criticism of politicians in general and personal reports of Covid-19 cure by hydroxychloroquine
<b>Post-public speech 5 (April 17)</b>	
1	Jokes about the return to the workplace
2	Criticism of the economic cost of the quarantine
3	Resignation of the health minister (Henrique Mandetta)
4	Jokes and criticism of the inauguration ceremony of the new health minister (Nelson Teich)
5	Support of treatments for Covid-19 including hydroxychloroquine and nitazoxanide

Next, we use the topic correlation via the overlap metric on these topics to calculate a similarity between them and to analyze how the rumors behave against the social discussions. For this analysis, we calculate the correlation between speech and social topics relative to the same speech.

Table 6 shows the correlation results of the 5 periods, relative to the five official speeches, again using the heat map for a better visualization. Results are truncated with a precision of two. It should be noted that in this case the cell colors of the heat map are such that the green indicates the more correlated or most similar while the red indicates the most distant in comparison with the others. This is the contrary of the Zika epidemic scenario because here the correlation is represented by the overlap coefficient while the first was represented by the KL-Divergence. Different shades of the colors were used here to facilitate the discerning of this important difference.

<b>Public Speech 1 Topics (March 12) X Post-Public Speech 1 Social Topics (March 13)</b>										
Speech Topic	Social Topic									
	1	2	3	4	5	6	7	8	9	10
A	0.08	0.04	0.08	0.17	0.04	0.04	0.04	0.17	0.08	0.38
<b>Public Speech 2 Topics (March 24) X Post-Public Speech 2 Social Topics (March 25)</b>										
Speech Topic	Social Topic									
	1	2	3	4	5					
A	0.31	0.46	0.00	0.15	0.31					
B	0.29	0.00	0.29	0.00	0.57					
C	0.22	0.00	0.44	0.00	0.00					
D	0.00	0.00	0.00	0.18	0.00					
E	0.00	0.00	0.00	0.00	0.00					
F	0.32	0.11	0.11	0.32	0.21					
G	0.17	0.50	0.33	0.00	0.33					

H	0.20	0.10	0.00	0.10	0.60		
I	0.50	0.25	0.75	0.50	0.50		
J	0.00	0.00	0.00	0.00	0.11		
<b>Public Speech 3 Topics (March 31) X Post-Public Speech 3 Social Topics (April 01)</b>							
Speech Topic	Social Topic						
	1	2	3	4			
A	0.13	0.13	0.13	0.25			
B	0.33	0.22	0.22	0.22			
C	0.22	0.00	0.33	0.44			
D	0.27	0.16	0.49	0.11			
E	0.14	0.00	0.14	0.43			
F	0.14	0.29	0.57	0.14			
G	0.25	0.50	0.25	0.00			
H	0.67	0.00	0.00	0.22			
I	0.06	0.06	0.52	0.39			
J	0.00	0.00	0.00	0.36			
K	0.31	0.00	0.77	0.62			
L	0.21	0.21	0.28	0.07			
<b>Public Speech 4 Topics (April 08) X Post-Public Speech 4 Social Topics (April 09)</b>							
Speech Topic	Social Topic						
	1	2	3	4	5	6	7
A	0.12	0.30	0.18	0.00	0.06	0.12	0.18
B	0.29	0.43	0.14	0.14	0.14	0.29	0.14
C	0.15	0.22	0.15	0.00	0.22	0.30	0.37
D	0.21	0.24	0.21	0.17	0.34	0.24	0.07
E	0.13	0.10	0.13	0.10	0.20	0.23	0.07
F	0.18	0.18	0.18	0.45	0.18	0.36	0.18
<b>Public Speech 5 Topics (April 16) X Post-Public Speech 5 Social Topics (April 17)</b>							
Speech Topic	Social Topic						
	1	2	3	4	5		
A	0.00	0.00	0.15	0.92	0.00		



B	0.20	0.12	0.16	0.41	0.16
C	0.13	0.10	0.13	0.17	0.10
D	0.09	0.13	0.17	0.17	0.26
E	0.00	0.14	0.14	0.35	0.07

### 5.3. Discussion

A throughout analysis of the results is presented here divided by the two scenarios used in the experiments, i.e., the Zika epidemic and the Brazilian political speeches.

#### 5.3.1. Zika Epidemic

For an overview of the main topics detected, we can cite at the first period of social topics, topic 1 that is related to a rumor of a possible relationship between a company (Monsanto) and microcephaly. Topic 2 is about the epidemic affecting the Olympic Games preparation to be held in Brazil. Topics 3, 4, and 7 which are related to the WHO declaration of the epidemic as a Public Health Emergency of International Concern, the news about travel warnings for pregnant women and the cases related in various countries as the disease spread respectively. Topics 5 and 8 are related to case discussions in other languages (Spanish and Portuguese specifically) on the most affected countries.

In the second period of the social topics, we have topic 5, which contains posts reporting the spread of the virus to other countries not initially affected. Topics 3 and 7 are related to measures to prevent the contamination and dissemination of the disease. Topic 9 is a local discussion of Zika in Brazil and topic 1 has the main information about the disease. Finally, topic 4 shows population concerns about the Matthew Hurricane that hit Central and North America in that period.

For the science topics, topics 1 and 5 suggests that the academic community was aware of the initial stages of the epidemic as the labels are associated with infection vectors, how it is transmitted, and the virus threat. Topic 2 refers to the time when the researchers already considered the disease an outbreak. Topics 3 and 4 refer to studies about the first cases in Brazil (Particularly Bahia, a state where local researchers identified the Zika virus for the first time in the region) and the relationship between Zika virus (ZIKV) and Dengue virus (DENV), another virus that shares some of the same transmission vectors.

In the second period of the scientific topics, we have discussions focused on prevention in topic 1, new forms of transmission like saliva in topic 7, and human semen in topic 4. Topic 2 shows the discovery of the association between the occurrence of microcephaly and Zika virus infection in pregnant women. Topic 6 suggests research between Zika and neurological problems caused by it like the Guillain–Barré syndrome.

Since the similar topics suggest a correlation between the social topics and the science topics, which is a reliable source, we will focus our discussion of rumors evidence in the most distant topics. Based on the average values of the heat map, we

analyze the ones with the highest average distance over all topics (values above 3) while still discussing particular cases when necessary.

Starting at the first period, the topic 1 shows a rumor that says that a Brazilian company (Monsanto) pesticides were the true cause of microcephaly in children, a disease associated with the Zika virus in many kinds of research and with clinical proof. This topic presented one of the highest distances overall (in relation to science topic 3). The topic 3 has a high distance related to science topic 4 too, but it seems like a dissimilarity issue because the first is related to the outbreak declaration and the latter with the relationship between the Dengue virus and the Zika virus.

Topic 4 has a high average and its related to local concerns about the virus spread reaching Miami and concerns about the Mathew hurricane, but it seems that this high average is caused by the topic being a local discussion of the Miami city and thus, unrelated to any science topic. This seems to be the case with the topic 5 and 8 too. Both have high averages and are linked to local discussions (Spanish and Portuguese speaking communities). Also, a curious case can be seen in topic 7 which has a low average but it is related to a car brand named “Zica” from the Tata company in India. The high average could be explained by the hashtags and the promoted terms used like “#zika”, “#zika virus”.

Finally, in the second period, we have Topic 1 that is dissimilar with the scientific topics concerning travels and neurologic diseases, but it has a relation with a rumor of vaccines causing babies to be infected with the disease.

Topics 6 and 8 have the highest average distances. Topic 6 is related to the rumor that a negative mutation in mosquitoes was causing the zika spread in the Americas. Topic 8 is a local community discussion in Brazil made in Portuguese.

The lowest averages are seen in topics 5 and 7. The first contains a mix or reliable information of the virus spreading to various countries along with some rumors of cases in countries that were not proven. The latter is associated with repellents marketing and relate to rumors of specific brands being more effective in combating the Aedes mosquito.

We present a summary of rumors found in the topics and rumors found in the reference rumor list [WHO 2016] to facilitate the discussion comprehension in Table 4.

**Table 4. Comparison of WHO rumors and rumor topics.**

Rumor Topic	WHO Rumor
Monsanto pesticides related to microcephaly (Period 1, Topic 1)	Evidence that pyriproxyfen insecticide causes microcephaly
Tata car named “Zica” that generated confusion with the virus name (Period 1, Topic 7)	Vaccines cause microcephaly in babies
Rumors of first cases in various countries (Period 2, Topic 5)	Most symptoms of Zika virus disease are equal from those of seasonal flu
Negative mutation in mosquitoes causing the spread of the virus (Period 2, Topic 6)	Bacteria used to control the male mosquito population are spreading Zika further
Repellent information and propaganda (Period 2, Topic 6)	Some repellents work better against the Aedes mosquito
	Evidence that Zika virus and its complications are linked to releases of genetically modified mosquitoes in

	Brazil
	Evidence that sterilized male mosquitoes contribute to the spread of Zika

The results of the analysis show that there were 5 topics related to rumors, with 3 of them related to rumors in the reference list and just one rumor in the list is not related with any topic (“Most symptoms of Zika virus disease are equal from those of seasonal flu”). Most of the rumor topics found had the highest averages (above 3,00 in both scenarios) or lowest averages (1,64 in the first scenario and 2.53 in the second) overall. Although one can conclude that the two extremes of the averages are evidence of false information there were also local discussion in other languages or limited to certain areas that had high averages too. Possibly this was caused by the different languages and thus different terms causing an increased distance over the topics.

### 5.3.2. Brazilian Political Speeches

Since there is a high number of topics over the five speeches, we summarize the main topics and findings relative to each speech to facilitate comprehension.

#### First speech - March 12th

The Brazilian President's key Speech Subject was the pandemic awareness in his first official speech with an appeal to his supporters to consider the situation at hand and stop the demonstrations in his favor scheduled for March 15th (Speech topic A). There were no evidences of rumors in his speech at this time according to the specialists. Post-speech, the highest correlation was with social topic 10, which consists of incentives for preventive measures such as washing hands, avoiding crowds, preserving social distance and strict hygiene habits.

#### Second speech - March 24th

In the second speech, because of the claim that the media is attempting to provoke social hysteria, Speech topic C was correlated with misinformation, since Brazil was considered to be less affected by the virus than the rest of the world. The president also argued that the country has a warm climate and that for healthy people, the Covid-19 symptoms are close to the normal flu.

Speech topic C was correlated with Social topics 1 and 3 which included people that minimized the pandemic's severity and accusing politicians of “using” the pandemic for their interests. In particular, social topic 3 deals with a part of the speech in which the President referred to the disease as a "small flu" ("gripezinha" in Portuguese), a phrase that minimized the seriousness of infected people's consequences. There were criticisms of politicians spreading social panic along these lines (Speech topic E). Despite showing lower correlation values, speech topic E is also correlated to Social topics 1 and 3.

The minimization of risks to "healthy" people (Speech topic F) was connected with discussions minimizing the seriousness of the pandemic and condemning political leaders who were expected to take advantage of the crisis (Social topic 1) and defending the positions of the president (Social topic 4).

The risk reduction for people outside the elderly groups (Speech topic G) was more closely linked to social topic 2, which discussed jokes about Prince Charles'

infection and an excerpt from the presidential speech in which Bolsonaro said that when contaminated, the military and athletes have mild symptoms.

The efficacy of the hydroxychloroquine (Speech topic H) was most related to social topic 5, related to hydroxychloroquine government authorization. Finally, the supposed cure for the disease (Speech topic I) was present in most social topics, with peak correlations between discussions about the pandemic minimization after the pronouncement (Social topic 3).

### **Third speech - March 31st**

Disinformation in the third speech was mainly related to the dichotomy between the economy and health (Speech topic B) and the reiterated argument on the effectiveness of hydroxychloroquine (Speech topic F). The first topic maintained a correlation with all topics, especially with the social topic 1 which contains criticism of local authorities and social distancing. The latter had a high correlation with social topic 3 which contains discussions about the importation of supplies to produce hydroxychloroquine in the country, and the military aid in the pandemic

The highest correlation among all topics was between speech topic K (economic recovery) and social topic 3, that is about the importation of supplies to produce hydroxychloroquine in the country by the Brazilian army.

### **Fourth speech - April 8th**

The fourth speech has three topics associated with disinformation, which we treat as rumors in this work. The first is on criticism of governors and their social distancing policy (Speech topic B). Being correlated particularly to social topic 2, which contains criticism of media coverage and measures for containing the spread of the virus.

The second rumor topic is about the dichotomy between the economy and health (Speech topic C). It is mostly correlated with Social topic 6, containing criticism of the WHO and the Minister of Health (Henrique Mandetta), and social topic 7, showing criticism of politicians, mainly of those in opposition parties.

The third and last rumor topic is also about the hydroxychloroquine efficacy (Speech topic D). Even though it has correlations with the majority of topics, social topic 5 has the highest correlation, with people defending the use of hydroxychloroquine and criticizing politicians who oppose the use of the drug.

### **Fifth Speech - April 16th**

Finally, the fifth speech has only one rumor related topic, which is again about the dichotomy between social distancing and economical recovery (Speech topic D). It was correlated mainly with social topic 4, about treatments, social distancing and the economic crisis.

Overall, it seems that the topics considered related or containing rumors by human specialists were more correlated with topics showing similar thematic. Unlike the previous experiment where our reference set contained authoritative data unrelated to rumors, in this case our reference data was annotated and described as rumors or not. In the Zika epidemic scenario the average correlation suggested the presence of rumors in the social discussions. In the scenario where we have the rumors topics described

beforehand, the average was not necessary and the simply correlation with the rumor topics already showed the social topics more related of propense to contain rumors.

Some other studies are necessary to verify that the patterns found here can be generalized in other cases and to study the types of rumors or fake information found in the data. It seems that the topics more unrelated to authoritative sources (like our science database) correspond to local discussions or rumors, while topics that are related to anything could be a deliberate fake or propaganda when used in marketing like repellent sellers or cars. Topics more related to authoritative sources that indicate the rumors (like the Presidential speeches) suggest the presence of terms of thematic related to rumors, possibly fake news and misleading social claims.

## 6. Conclusions

The detection of false information had always been a topic that attracted much attention over the years. The modern threats of social networks, “fake news”, and the spread of viral false or unverified information demand different methods to detect rumor and analyze them to cope with the dynamic nature of the online communities.

Most of the works in the area have been trying to detect false messages but few of them were concerned with a topic level detection, i. e., which subjects are related to rumors. This work proposes a topic model approach to detect these subjects and uses topic correlation via KL-Distance and Overlap coefficient between two domains: A social network domain and a “ground truth” domain where the topics are verified (in our case a scientific domain and a political domain).

We conduct a study of how the relationship between topics of these different domains relate to each other in the Zika epidemic scenario and the Brazilian presidential speeches scenario. In the first, we analyze how they relate to verified rumors from the World Health Organization. The last was compared with rumors present in the speeches itself, according to communication specialists. We compare the results found and our findings suggest evidence that topics that have very low or high correlation could indicate rumors or limited local discussions (most in other languages). Our main contributions are thus: 1) The cross-topic model method for inferencing rumors; 2) The topic correlation approach to detect rumor topics; and 3) The analysis of the metric differences between localized discussions, rumors, and the behavior of rumor topics.

This work open opportunities in the detection of fake news, disinformation or topics and labels that are indicative of unreliable information that could be used to assess message reliability, the impact of a rumor in the most important subjects of the social network and the very detection of rumor topics that are widespread in the social networks.

## References

- Ahsan, M. and Kumari, M. (2019). Rumors and their controlling mechanisms in online social networks: A survey. *Online Social Networks and Media* [[GS Search](#)]
- Al-Khalifa, H. S. and Al-Eidan, R. M. (2011). An experimental system for measuring the credibility of news content in Twitter. *International Journal of Web Information Systems*, v. 7, n. 2, p. 130–151. [[GS Search](#)]
- Allport, G. and Postman, L. (1947). The psychology of rumor. [[GS Search](#)]

- Blei, D., Carin, L. and Dunson, D. (2010). Probabilistic topic models. *IEEE Signal Processing Magazine*, v. 27, p. 55–65. [[GS Search](#)]
- Blei, D. and Lafferty, J. (2009). Topic models. : *classification, clustering, and applications*, [[GS Search](#)]
- Blei, D. M., Ng, A. Y. and Jordan, M. I. (2003). Latent Dirichlet Allocation. *Journal of Machine Learning Research*, v. 3, n. 4–5, p. 993–1022. [[GS Search](#)]
- BrasilGov, T. (2020a). Pronunciamento oficial do Presidente da República, Jair Bolsonaro (12/03/2020) - YouTube. <https://www.youtube.com/watch?v=bS2qiXHtMnI>, [accessed on Feb 15].
- BrasilGov, T. (2020b). Pronunciamento do presidente da República, Jair Bolsonaro (24/03/2020) - YouTube. [https://www.youtube.com/watch?v=VI\\_DYb-XaAE](https://www.youtube.com/watch?v=VI_DYb-XaAE), [accessed on Feb 15].
- BrasilGov, T. (2020c). Pronunciamento do presidente da República, Jair Bolsonaro (31/03/2020) - YouTube. [https://www.youtube.com/watch?v=16RR2rG\\_AKA](https://www.youtube.com/watch?v=16RR2rG_AKA), [accessed on Feb 15].
- BrasilGov, T. (2020d). Pronunciamento do presidente da República, Jair Bolsonaro (08/04/2020) - YouTube. <https://www.youtube.com/watch?v=x04OKkxT2Tc>, [accessed on Feb 15].
- BrasilGov, T. (2020e). Presidente Jair Bolsonaro faz pronunciamento (16/04/2020) - YouTube. <https://www.youtube.com/watch?v=GwiVPFZ5610>, [accessed on Feb 15].
- Cao, J., Guo, J., Li, X., et al. (10 jul 2018). Automatic Rumor Detection on Microblogs: A Survey. [[GS Search](#)]
- Castillo, C., Mendoza, M. and Poblete, B. (2011). Information credibility on Twitter. In *Proceedings of the 20th International Conference Companion on World Wide Web, WWW 2011*. [[GS Search](#)]
- DiFonzo, N. and Bordia, P. (2007). *Rumor psychology: Social and organizational approaches*. [[GS Search](#)]
- Fundação Oswaldo Cruz., R. and Araujo, I. S. (2016). *A mídia em meio às 'emergências' do vírus Zika: questões para o campo da comunicação e saúde*. Fundação Oswaldo Cruz. v. 10 [[GS Search](#)]
- Kullback, S. and Leibler, R. A. (1951). On Information and Sufficiency. *The Annals of Mathematical Statistics*, v. 22, n. 1, p. 79–86. [[GS Search](#)]
- Kwon, S., Cha, M., Jung, K., Chen, W. and Wang, Y. (2013). Prominent features of rumor propagation in online social media. In *Proceedings - IEEE International Conference on Data Mining, ICDM*. [[GS Search](#)]
- Mendoza, M., Poblete, B. and Castillo, C. (2010). Twitter under crisis: Can we trust what we RT? In *SOMA 2010 - Proceedings of the 1st Workshop on Social Media Analytics*. [[GS Search](#)]
- Nolasco, D. and Oliveira, J. (2016). Detecting knowledge innovation through automatic topic labeling on scholar data. In *Proceedings of the Annual Hawaii International Conference on System Sciences*. [[GS Search](#)]

- Nolasco, D. and Oliveira, J. (2018). Subevents detection through topic modeling in social media posts. *Future Generation Computer Systems*, [GS Search]
- Nolasco, D. and Oliveira, J. (2020). Mining social influence in science and vice-versa: A topic correlation approach. *International Journal of Information Management*, v. 51. [GS Search]
- Peterson, W. A. and Gist, N. P. (sep 1951). Rumor and Public Opinion. *American Journal of Sociology*, v. 57, n. 2, p. 159–167. [GS Search]
- Ratkiewicz, J., Meiss, M., Conover, M., et al. (2011). Detecting and Tracking Political Abuse in Social Media. In *Proceedings of the Fifth International AAAI Conference on Weblogs and Social Media*. [GS Search]
- Sharma, K., Qian, F., Jiang, H., et al. (2019). Combating fake news: A survey on identification and mitigation techniques. *ACM Transactions on Intelligent Systems and Technology* [GS Search]
- Steyvers, M. and Griffiths, T. (2007). Probabilistic topic models. *Handbook of latent semantic analysis*, [GS Search]
- Takahashi, B., Tandoc, E. C. and Carmichael, C. (2015). Communicating on Twitter during a disaster: An analysis of tweets during Typhoon Haiyan in the Philippines. *Computers in Human Behavior*, v. 50, n. 2015, p. 392–398. [GS Search]
- Vijaymeena, M. . and Kavitha, K. (2016). A Survey on Similarity Measures in Text Mining. *Machine Learning and Applications: An International Journal*, v. 3, n. 1, p. 19–28. [GS Search]
- WHO (2016). WHO - Dispelling rumours around Zika and complications. <http://www.who.int/emergencies/zika-virus/articles/rumours/en/>, [accessed on Apr 29].
- WHO (2020). COVID-19 advice - Mythbusters | WHO Western Pacific. <https://www.who.int/westernpacific/emergencies/covid-19/information/mythbusters>, [accessed on Feb 15].
- Zubiaga, A., Aker, A., Bontcheva, K., Liakata, M. and Procter, R. (1 feb 2018). Detection and resolution of rumours in social media: A survey. *ACM Computing Surveys*. Association for Computing Machinery. [GS Search]