

**POLARIZAÇÃO POLÍTICA E O IMPEACHMENT  
DE 2016: UMA ANÁLISE DE DADOS REAIS E DE  
MÍDIAS SOCIAIS**

ROBERTA COELI NEVES MOREIRA

**POLARIZAÇÃO POLÍTICA E O IMPEACHMENT  
DE 2016: UMA ANÁLISE DE DADOS REAIS E DE  
MÍDIAS SOCIAIS**

Dissertação apresentada ao Programa de Pós-Graduação em Ciência da Computação do Instituto de Ciências Exatas da Universidade Federal de Minas Gerais como requisito parcial para a obtenção do grau de Mestre em Ciência da Computação.

ORIENTADOR: GISELE LOBO PAPP  
COORIENTADOR: PEDRO OLMO STANCIOLI VAZ DE MELO

Belo Horizonte  
Dezembro de 2018

ROBERTA COELI NEVES MOREIRA

**POLITICAL POLARIZATION AND THE  
IMPEACHMENT OF 2016: AN ANALYSIS OF  
REAL-WORLD AND SOCIAL-MEDIA DATA**

Dissertation presented to the Graduate Program in Computer Science of the Federal University of Minas Gerais in partial fulfillment of the requirements for the degree of Master in Computer Science.

ADVISOR: GISELE LOBO PAPPA  
CO-ADVISOR: PEDRO OLMO STANCIOLI VAZ DE MELO

Belo Horizonte

December 2018

© 2018, Roberta Coeli Neves Moreira  
Todos os direitos reservados

Ficha catalográfica elaborada pela Biblioteca do ICEx - UFMG

Moreira, Roberta Coeli Neves.

M837p Political polarization and the impeachment of 2016: an analysis of real-world and social-media data / Roberta Coeli Neves Moreira — Belo Horizonte, 2018.  
xxiv, 66. il.; 29 cm.

Dissertação (mestrado) - Universidade Federal de Minas Gerais – Departamento de Ciência da Computação.

Orientadora: Gisele Lobo Pappa  
Coorientador: Pedro Olmo Stancioli Vaz de Melo

1. Computação – Teses. 2. Polarização política. 3. Redes Sociais. 4. Análise de sentimentos. I. Orientadora. II. Coorientador III. Título.

CDU 519.6\*82.10 (043)



UNIVERSIDADE FEDERAL DE MINAS GERAIS  
INSTITUTO DE CIÊNCIAS EXATAS  
PROGRAMA DE PÓS-GRADUAÇÃO EM CIÊNCIA DA COMPUTAÇÃO

## FOLHA DE APROVAÇÃO


Political Polarization and the Impeachment of 2016: an analysis of real-world  
and social-media data

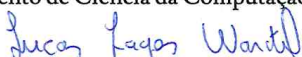
**ROBERTA COELI NEVES MOREIRA**

Dissertação defendida e aprovada pela banca examinadora constituída pelos Senhores:

  
PROFA. GISELE LOBO PAPPÁ - Orientadora  
Departamento de Ciência da Computação - UFMG

  
PROF. PEDRO OLMO STANCIOLI VAZ DE MELO - Coorientador  
Departamento de Ciência da Computação - UFMG

  
PROFA. MIRELLA MOURA MORO  
Departamento de Ciência da Computação - UFMG

  
PROF. LUCAS LAGES WARDIL  
Departamento de Física - UFMG

Belo Horizonte, 22 de fevereiro de 2019.

# Agradecimentos

Durante o período do mestrado, a vida trouxe muitas mudanças e desafios. Por isso, muito tenho a agradecer a tantos à minha volta pela finalização desta etapa.

Primeiramente, agradeço aos meus pais pelo incentivo que sempre me deram nos estudos. Sou eternamente grata pelas tantas abdições em nome da minha educação e pelo apoio emocional nessa trajetória. Ao Alain, agradeço muito pela compreensão nos altos e baixos desse período, pelo apoio e carinho nos instantes de desânimo. Aos familiares queridos, os quais também foram um suporte fundamental nessa estrada, agradeço pelos momentos de descontração e alegria.

Ao Gabriel, meu amigo-irmão, com quem tive a oportunidade de dividir mais essa etapa, agradeço pela amizade sincera, pelas parcerias nos trabalhos e pelas dificuldades compartilhadas. À Malu, quem o mestrado me deu a chance de conhecer melhor, agradeço pela companhia nos trabalhos e estudos e, claro, pelas conversas das quais tanto tenho saudade. A todos os meus amigos, especialmente à Thaís e à Marina, pelas boas energias e pela compreensão em meus momentos de ausência.

Aos meus orientadores, Gisele e Pedro, agradeço pela oportunidade e pela imensa paciência. Muito obrigada pela compreensão com as tantas mudanças nesse período e por sempre estarem dispostos a ajudar, mesmo à distância. Agradeço também pelo incentivo e pelo apoio diante das dificuldades que encontrei ao fazer a pesquisa, bem como pelas ideias e sugestões valiosas.

A todos os professores, agradeço pelo trabalho árduo de fazer ciência no Brasil e, ao mesmo tempo, pela dedicação em nos proporcionar uma educação de qualidade. A todos os funcionários do DCC, que contribuem para o bom funcionamento de toda a estrutura acadêmica, muito obrigada!

Agradeço também à Google pelas oportunidades e pelo auxílio financeiro durante o mestrado.

*“Acautelai-vos contra os juízos arrebatados pela paixão porque esta desfigura muitas vezes a verdade. Aquele que olha por um vidro de cor vê todos os objetos da cor desse vidro: se o vidro é vermelho, tudo lhe parece rubro; se é amarelo, tudo lhe apresenta completamente amarelado. A paixão está para nós como a cor do vidro para os olhos.”*

(Malba Tahan em “O homem que calculava”)

# Resumo

Assuntos políticos frequentemente geram acalorados debates ao redor do mundo, os quais tendem a aumentar a divergência de opinião dos indivíduos. A emergência de tais ideias opostas descreve a polarização política, fenômeno que tem sido evidenciado e catalisado pela massificação da Internet e das redes sociais nos últimos anos. No Brasil, a polarização política se intensificou com o processo de impeachment de Dilma Rousseff e seus eventos atrelados, cujo duelo de ideias entre os cidadãos e entre as elites partidárias pôde ser visto tanto nas ruas quanto nas mídias sociais. Tendo em vista a crescente influência de tais mídias na formação de opinião dos cidadãos, a presente pesquisa visou criar um arcabouço teórico e ferramental para análise de dados reais e de mídias sociais para estudo da polarização política, com foco no processo de impeachment ocorrido no Brasil em 2016. Nossa metodologia envolveu 4 principais partes: I) a análise temporal de tópicos a partir de tweets de parlamentares brasileiros; II) a quantificação das polaridades e da polarização dos cidadãos brasileiros no Twitter e dos deputados federais em suas votações na Câmara; III) a correlação entre a polarização dos deputados e seus tópicos no Twitter; e IV) o estudo da associação entre a polarização dos deputados e a polarização do público. Nossos resultados mostraram que a crise política e as atividades dos parlamentares foram os tópicos de maior relevância dentre os discutidos pelos políticos no Twitter. Quanto à polarização, os deputados apresentaram um aumento da polarização após dezembro de 2015, o que coincidiu com o início do processo de impeachment na Câmara. O público geral registrou altos valores de polarização durante todo o período de 2016, apresentando uma polarização mais alta que os deputados em todos os meses. Uma análise das variáveis relacionadas à métrica de polarização nos permitiu também observar que a polarização dos políticos foi mais influenciada pelo nível de divergência entre as opiniões centrais dos grupos de ideologias opostas. Por outro lado, as pequenas variações nos valores da polarização do público se deveram à diferença no tamanho dos grupos de opiniões contrárias, mostrando que o fenômeno de polarização entre os brasileiros está mais relacionado à quantidade de indivíduos que se juntam a cada dos grupos opostos.



# Abstract

Political events are often topics of heated discussions around the globe, which tend to increase the opinion divergences among individuals. The emergence of these contrasting ideas defines the political polarization, a phenomenon which have been evidenced and boosted by the popularization of Internet access and social networks over the last few years. In Brazil, the political polarization was intensified by the impeachment of Dilma Rousseff and its related events, revealing ideological conflicts among the citizens and the political elites not only in the streets but also in online social media. In view of the growing influence of the social platforms in the opinion formation of people, our work aimed at developing computational methods to analyze real-world and social-media data in order to study political polarization, with a focus on the impeachment proceedings of 2016 in Brazil. Our methodology involved 4 main parts: I) the temporal topic evolution from the tweets posted by Brazilian representatives; II) the measurement of the polarities and the overall polarization of the Brazilian general public in Twitter and of the representatives in their votes on bills at the Lower House; III) the correlations between the polarization of the politicians and their discussed topics on Twitter; and IV) an analysis of the associations between the politicians and the general public polarization. Our results showed that the political crisis and the activities of the representatives were the most relevant topics discussed by the politicians on Twitter. Regarding the polarization analysis, we observed that the politicians polarization increased after December of 2015, coinciding with the launch of the impeachment proceedings at the Lower House. The general public presented high values of polarization during the whole period, recording higher values when compared to the representatives. By analysing the variables related to the polarization metric, we also observed that the polarization among politicians was more affected by the level of divergence between the central opinions of the main opposite groups. On the other hand, small fluctuations in the people polarization are related to differences in the size of the contrary groups, meaning that the polarization phenomenon among Brazilian citizens is more influenced by the number of individuals which join any of the ideologically opposed groups.

# List of Figures

1.1	Timeline of the main events regarding the impeachment of Dilma Rousseff.	3
3.1	Example of probability density distribution of polarities, showing the populations of opposite opinions ( $A^-$ and $A^+$ ), the gravity centers of each population ( $gc^-$ and $gc^+$ ) and the distance between these gravity centers ( $d$ ).	17
3.2	Example of probability density function having maximum polarization index ( $\mu = 1$ ).	20
3.3	Example of probability density function having a population that is not polarized ( $\mu = 0$ ).	20
3.4	Example of probability density function where $d = 1$ and $\Delta A > 0$ .	21
3.5	Example of probability density function where $0 < d < 1$ and $0 < \Delta A < 1$ .	22
4.1	Graphical Representation of Biterm Topic Model (BTM).	25
4.2	Example of application of the proposed topic evolution method.	26
4.3	Number of tweets over time for the politicians.	29
4.4	Participation of the politicians on Twitter according to the political spectrum of their parties ( $P_s$ ) over time.	30
4.5	Word clouds for monthly datasets.	32
4.6	Relevance of super topics $T'_p$ over the months.	32
4.7	Word Relevance over the months - Super Topic 1 (Political Crisis).	34
4.8	Overview of the steps to calculate politicians polarization.	36
4.9	Average polarity of the Brazilian parties over time.	38
4.10	Probability density functions for the polarity values of politicians per month.	39
4.11	Time evolution of polarization index (c), and its related variables: difference in population sizes (a) and distance between gravity centers (b).	41
4.12	Time Evolution of Politicians Polarization Index and Percentage of Tweets of a Topic.	43
5.1	Overview of the steps to calculate the polarities of the general public.	47

5.2	Probability density functions for polarity values of Brazilian general public per month. . . . .	50
5.3	Time evolution of polarization index $\mu$ and its related variables for general public analysis. . . . .	51

# List of Tables

- 4.1 Words describing the super-topics for politicians. . . . . 33
- 4.2 Spearman Correlation Coefficient between Politicians' Percentage of Topics and Polarization Index. . . . . 44
  
- 5.1 Keywords used to collect tweets for the general public. . . . . 46
- 5.2 Hashtags related to the impeachment event in the general public dataset. . 48
- 5.3 Statistical summary for the difference between populations sizes ( $\Delta A$ ), distance between gravity centers ( $d$ ) and polarity index ( $\mu$ ) in people and politicians studies. . . . . 54
- 5.4 Comparison between polarization index  $\mu$ , difference between populations sizes  $\Delta A$  and distance between gravity centers  $d$  for people and politicians. 54

# Contents

Agradecimientos	v
Resumo	vii
Abstract	viii
List of Figures	ix
List of Tables	xi
<b>1 Introduction</b>	<b>1</b>
1.1 Contextualization . . . . .	3
1.2 Motivation . . . . .	5
1.3 Objectives . . . . .	6
1.4 Organization of Dissertation . . . . .	7
<b>2 Related Work</b>	<b>8</b>
2.1 Topic Evolution Over Time . . . . .	8
2.2 Polarity Evaluation . . . . .	10
2.3 Polarization Measurement . . . . .	12
2.4 Computational Social Science and Politics . . . . .	13
<b>3 Polarization Index</b>	<b>16</b>
3.1 Background . . . . .	17
3.2 Difference between populations sizes . . . . .	18
3.3 Distance between gravity centers . . . . .	18
3.4 Polarization Index . . . . .	19
<b>4 Polarization of Politicians</b>	<b>23</b>
4.1 Investigating Topic Evolution in Social Media . . . . .	23

4.1.1	Data . . . . .	24
4.1.2	Methods . . . . .	24
4.1.3	Experimental Results . . . . .	29
4.2	Calculating Polarization from Roll-Call Voting Data . . . . .	34
4.2.1	Data . . . . .	35
4.2.2	Methods . . . . .	35
4.2.3	Experimental Results . . . . .	37
4.3	Correlations . . . . .	42
<b>5</b>	<b>Polarization of People</b>	<b>45</b>
5.1	Calculating Polarization in Social Media . . . . .	45
5.1.1	Data . . . . .	46
5.1.2	Methods . . . . .	47
5.1.3	Experimental Results . . . . .	50
5.2	Comparing People and Politicians Polarization . . . . .	52
<b>6</b>	<b>Conclusions and Future Work</b>	<b>55</b>
6.1	Future Works . . . . .	57
	<b>Bibliography</b>	<b>59</b>

# Chapter 1

## Introduction

Religion, sports and politics are often topics of heated discussions around the globe. These debates often reveal contrasting ideas about the discussed matter and, in this situation, when people see what others think, they tend to strengthen their prior beliefs [Sunstein, 2002]. People that are contrary to abortion legalization, for example, tend to be more extremely opposed to it after interacting with others that share the same point of view. On the other hand, supporters of the cause reinforce their position after communicating with other pro-legalization individuals, which intensifies the controversy around the topic. In Social Sciences, this simultaneous presence of conflicting tendencies or principles characterizes the process of group polarization [Fiorina and Abrams, 2008]. Above the many subjects that raises polarization in society, politics has been shown to be one of the most fertile grounds to bring disagreements between people into the open.

In the field of politics, social scientists describe two types of polarization: elite polarization and mass polarization. Elite polarization is characterized by the high ideological discrepancy between political parties and the strong similarity of positions inside parties [Druckman et al., 2013]. Mass polarization, in turn, is related to the segregation of common individuals in society due to the divergence of opinions regarding actions and ideas of the political elites [Baldassarri and Gelman, 2008]. With respect to this process of social division, the American society and political system are one of the most widely studied examples [Fiorina and Abrams, 2008; Abramowitz and Saunders, 2008; DiMaggio et al., 1996]. Besides the intensification of the ideological conflicts between Democrats and Republicans – the two main political parties in the United States – Abramowitz and Saunders [2008] explain that mass polarization has also grown in the country, backed up by an increase in the education level of the population and access to information, especially with the advent of the Web.

The Web allows political parties and the population to spread their opinions quickly and to a large audience, evidencing and amplifying the political polarization process around the world [Adamic and Glance, 2005; Farrell and Drezner, 2008; Farrell, 2012]. Regarded as one of the main factors that contributed to this process online, social networks do not only act as a vehicle to consume information, but also make possible for people to express their positions and participate on political campaigns and manifestations. By doing so, they reinforce previous beliefs of people by connecting like-minded individuals, but also set the stage for energetic arguments between the citizens that have opposite points of view. In Brazil, events such as the protests of 2013, the elections of 2014 and the impeachment proceedings of Dilma Rousseff mobilized people and divided opinions. Ruediger et al. [2014] showed that these disputes have emerged due to the rise of Internet usage over the last few years in the country. According to the study, online social media platforms have been valuable tools to demand improvements on public services, contest representatives about their actions and to integrate people having common interests, which were essential ingredients for the outbreak of protests on June of 2013.

Online social networks disseminate texts, images and videos from many different sources, gathering a large amount of information in a single place. For this reason, the popularization of this sort of media plays an important role in the opinion formation of the citizens, especially with regard to politics. The influence of such media has been endorsed by recent statistics from IBOPE [2016]: in Brazil, approximately 51% of the voters consume political information from Facebook, Twitter or Whatsapp. Among these, 27% of people stated that they had a more favorable impression of a politician or party after viewing posts on online social networks. By comparison, 56% of the individuals declared that they changed their opinion for the worse about politicians and political parties due to what they read on these platforms.

Considering their high impact over people, it is essential to understand the dynamics of online social networks since they allow the emergence of the so-called “echo chambers”. In the context of social media, an echo chamber can be defined as an environment in which the user only reads information and discusses ideas that coincide with her own opinions and interests, ignoring or blocking posts and other users that share alternative views. According to Garrett [2009], these echo chambers arise because people are more interested in content that supports their points of view instead of content that challenges their beliefs. The author explains that the large amount of information and resources provided by the Web and online social media are increasing the ideological fragmentation among individuals by creating these “filter bubbles”.

In view of the power of such technologies on shaping the political beliefs of the



citizens, sociologists and computer scientists see a great potential in online social networks for understanding society and predicting outcomes in politics [Farrell, 2012]. This is because, besides providing a massive amount of data in real-time, these platforms make it possible to reach a considerable part of the population in the research studies. In sociology, some interesting examples include the studies of Gerbaudo [2018] about the use of online social media by political activists as a part of a project of re-appropriation of public space; and the works of Wolfsfeld et al. [2013] about the role of these platforms in the collective actions of the Arab Spring protests. Besides sociological studies, there are a number of works in the computer science literature that analyze such media in order to classify, characterize and even predict people reactions regarding political events, exploring computational methods to deal with the huge volume of data generated by the dynamism of these online social networks.

## 1.1 Contextualization

This section introduces a summary of the context around the impeachment proceedings of Dilma Rousseff in 2016, which is our event of interest in this work. Figure 1.1 shows a timeline of the main episodes concerning the impeachment. The events shown in the timeline of the Figure 1.1 and described below are based on the works of Tatagiba [2018] and Velasco et al. [2016].

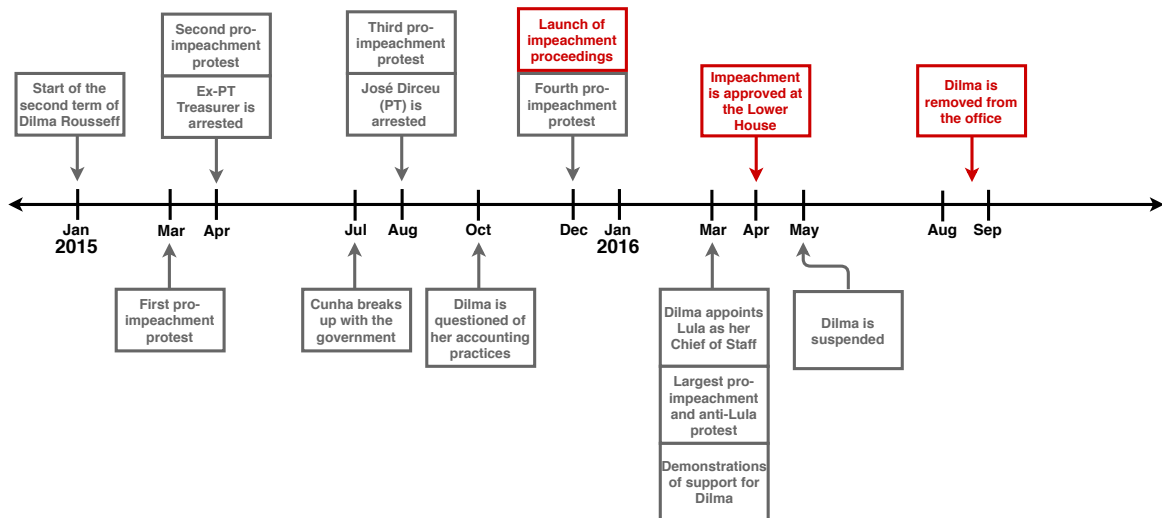


Figure 1.1: Timeline of the main events regarding the impeachment of Dilma Rousseff.

In her second term as president, which started in January of 2015, Dilma Rousseff had to deal with many challenges in the Brazilian social, economic and political

scenarios. After facing protests from the population on her previous term, she won a tight election race against Aécio Neves – candidate of Brazilian Social Democracy Party (PSDB), her major opposition party [Watts, 2014]. Aécio questioned election results and his party continuously contested the legitimacy of Dilma right after her reelection, which, together with the most fragmented congress in the history of the country, showed that the president would have to manage an unstable political system.

Besides this political crisis, right after Dilma started her second mandate, there were new revelations of politicians involved in the largest corruption scandal in the history of Brazil: “Lava Jato” (*Car Wash* operation). The operation consisted of a criminal investigation, carried out by the Federal Police of Brazil, to account for money laundering and corruption in Petrobras, the Brazilian largest state-owned oil company. The operation revealed a complex bribery scheme and put in jail big players of the Brazilian political scenario.

In the beginning of March of 2015, the General Attorney of the Republic revealed a list of politicians under suspicion of corruption in the *Car Wash* operation, which was sent to the Federal Supreme Court for further investigations. Among many of the investigated politicians, there were influential members of the Workers Party (PT), party to which Rousseff belongs to. In addition to the continuous hostilities instigated by Dilma’s opponents, these suspicions on her party members fueled even more tension and social discontent on Brazilians, contributing to undermine the prestige of the re-elected president. Under these circumstances, more than 2 million people took the streets to protest against corruption and to demand the impeachment of Dilma on March 15, 2015. Shortly after, on April 12 of the same year, another march brought thousands of people to the streets, motivated by the same reasons of the first one. Both demonstrations were organized through online social networks by political movements such as “Vem pra Rua” (*Come to the Streets*) and “Movimento Brasil Livre” (*Free Brazil Movement*) [Bedinelli and Martín, 2015]. A few days after the protests, João Vaccari Neto, PT treasurer, was arrested for his involvement in the Petrobras bribery scheme.

On July 17 of 2015, Eduardo Cunha, the Head of the Lower House of Congress, announced his break-up with the government and called himself an opponent of the government of Dilma. Not much later, at the beginning of August of 2015, José Dirceu, another important member of the Workers Party and ex-minister on previous mandates, was also arrested by the Federal Police under accusations of corruption. These sequence of arrests and the huge volume of corruption reports, combined with the questioning of the president accounting practices by the Federal Accounting Tribunal on October of 2015, deepened the political and economical instability of the government.

With the claim of the allegations of fiscal pedaling against Dilma, Hélio Bicudo

and other jurists filed an impeachment request to the House of Representatives on September 1, 2015. Also accused of corruption in the Petrobras scheme, Cunha was the one responsible for analysing and accepting these requests. On December 2, 2015, he agreed to launch the impeachment proceedings against the president. A few days later, thousands of Brazilian protesters took the streets again to demand the ouster of Dilma Rouseff.

The *Car Wash* investigations continued in 2016 and, besides the impeachment proceedings, Rouseff had to deal with new turbulences regarding her party. In addition to the arrest of her marketers over illicit campaign funding, Delcídio do Amaral – ex-leader of PT in Senate – testified that Rouseff and Lula tried to actively obstruct the Lava Jato investigations on corruption. On March 4, 2016, former president Lula, Dilma’s political mentor and closest ally, was questioned by the police about personal benefits from construction companies. A few days later, Rouseff invited Lula to join her cabinet, leading to accusations that she planned to nominate him so as to protect her ally from prosecution. These events sparked the largest anti-government protest in the history of the country: more than 3.5 million people reached the streets to demand the impeachment of the president. In this same period, a large number of protesters also gathered across the country to demonstrate support for Dilma, advocating against the impeachment.

On April 17, 2016, 367 out of 513 representatives in the Lower House of Congress voted to impeach the president. The process moved to the Senate, where the politicians started the impeachment proceedings and suspended Dilma from the office for 180 days. After she left, Michel Temer assumed as the interim president of Brazil. By August of the same year, Rouseff went to the Congress to present her defense, facing the senators to answer questions regarding the charges against her. After her defense delivered the final speeches, on August 31, the Senate voted to remove Dilma from the office with 61 votes in favor and 20 against impeachment. Temer was then officially inaugurated as president of Brazil and was in charge until the end of 2018.

## 1.2 Motivation

In Brazil, the impeachment proceedings of Dilma Rouseff in 2016 and its related events had a huge repercussion on the political, economic and social scenarios. Besides intensifying the disagreements among political parties, the sequence of events motivated common citizens to adopt an opinion concerning the removal of the president from office, which led a substantial number of Brazilians to the streets either to demand the

impeachment of Dilma or to protest against it. This conflict of ideas was also seen on online social networks, which reflected the points of view of the individuals and, at the same time, increased the differences between the opposite opinions and set the stage for heated political debates among users and parties [Ribeiro and Gomes Goveia, 2016].

In the view of the political polarization and further effects caused by the impeachment proceedings and other political events, it is important to understand the dynamics of how content is generated and propagated through online social networks, especially due to the growing influence of such platforms in the opinion formation of people. These studies have sociological relevance, since they contribute to a better understanding of society by delineating differences between the perspectives of the social groups, as well as by identifying the actual interests behind the opinions of these groups. Such analyses are also pertinent to discover strategies to reduce the impact of the “echo chambers” on social media, since it is necessary to diversify the ideas that are spread through these channels to preserve the integrity of the democratic institutions [Garrett, 2009; Garimella et al., 2016]. Furthermore, research in this area can help to identify the causes of segregation and political conflicts that emerge from ideological divergences in the virtual world, making it possible to find solutions to reduce online confrontations and avoid their transposition to the real world.

In the computer science area, these studies are also relevant, since they explore data mining and natural language processing methods in order to extract knowledge from a large volumes of data. In this context, they also face the challenge of the nature of the data: online social networks convey short and informal messages, which are difficult to process since they lack enough content to derive information from.

### 1.3 Objectives

Given the contextualization and the motivations, the main goal of this work is to develop computational methods to analyze online and offline data in order to study political polarization. Our focus here is to investigate real and virtual data regarding the impeachment proceedings of 2016 in Brazil, for which we tried to comprise a large number of individuals for a more thorough study.

The theoretical part of this study consisted of identifying and understanding the sociological and linguistical aspects which characterize the different opinions that are shared on Twitter with regard to the impeachment. The methodological part of our research involved developing approaches to deal with the computational challenges of

extracting information from massive amounts of data, particularly when it comes to short and informal texts from social media which processing is a difficult task.

Considering these aspects, the specific objectives of this study are:

- Analysis of the behavior of Brazilian representatives on the Web by developing methods to identify their main discussed topics on social media and how these topics changed over the time period around the impeachment event.
- Definition of metrics to evaluate political opinions expressed by the Brazilian general public on social media posts and by the Brazilian politicians through their actions on the real world, investigating online and offline data to find the main aspects that determine the position concerning the impeachment or its characters.
- Development of approaches to automatically identify and quantify the polarity of individuals holding the identified political opinions.
- Definition of metrics to quantify political polarization among the Brazilian general public and politicians, also investigating its temporal changes and identifying variables that are related to the phenomena.
- Analysis and application of methods to compare online and offline data, in order to investigate possible associations between them and understand the impact of social media over actions and behaviors in the real world.

## 1.4 Organization of Dissertation

This dissertation is organized as follows. Chapter 2 presents some of the works from the literature that are related to our research. Chapter 3 gives a more detailed description of the polarization model that we use to quantify group polarization in our study. Chapter 4 introduces our study of the polarization of the politicians, describing the analysis of their temporal topic evolution on social media and the evaluation of their group polarization over time. Chapter 5 presents our study of the polarization of the Brazilian people, explaining how we calculate the polarity of each individual in the dataset and the derived group polarization for the general public, as well as showing a comparison between the studies of the politicians and the people. Finally, Chapter 6 concludes this work and describes future directions.

# Chapter 2

## Related Work

Our study comprises several methods to evaluate the political polarization in virtual and real-world datasets. It involves detecting the topics discussed by the Brazilian politicians over time, as well as quantifying the individuals' polarity and the group polarization for both the politicians and the general public. For this reason, this chapter presents previous work on the four main sides of our research: topic evolution over time (Section 2.1), polarity evaluation (Section 2.2), polarization quantification (Section 2.3) and studies on computational social science and politics (Section 2.4).

### 2.1 Topic Evolution Over Time

The Web has become one of the main vehicles people use to consume information and its growing public is constantly reading, creating and sharing new content. Due to the large volume of textual data, detecting topics in these messages is important to understand the central subjects that are discussed in the virtual world. To achieve this, one can use topic modeling algorithms, which are statistical methods designed to discover topics in a collection of documents. In these probabilistic models, a topic is characterized by a probability distribution over a vocabulary, and a document consists of a mixture of these topics. One of the most widely used models for topic identification is the Latent Dirichlet Allocation (LDA), proposed by [Blei et al. \[2003\]](#), which is used for summarizing large text collections. However, with regard to topic modeling on social media, one of our interests in this work, there are other algorithms that explore the short nature of the texts and deal with the data sparsity problem that is not contemplated by the conventional topic models. To that end, [Yan et al. \[2013\]](#) presented Biterm Topic Model (BTM), which finds topics by directly modeling the generation of biterns – pairs of co-occurring words – in a corpus of short texts. BTM receives

as input the number of topics to be extracted and the set of biterns from the whole corpus, and produces as output the probability of each topic, as well as the probability of a word in the vocabulary given a topic.

Besides dealing with the large volume of texts by finding topics from a corpus of documents, as data is produced at a high speed, it is also important to consider that topics may change over time. For this reason, investigating topic modifications is relevant to study the causes of the rise and decline of subjects at a certain time period. In this context, some studies used document clustering to model topic evolution; other works developed variants of the traditional probabilistic topic models – such as LDA – adding time as an additional factor so as to find topic variations over time.

As one of the studies that adopts clustering techniques, [Mei and Zhai \[2005\]](#) applied a general probabilistic model to discover theme patterns and generated a graph to find word clusters for each time period. This graph – so-called evolutionary theme graph – is used to determine how topics change over time and how previous topics influence later ones. Other proposed models that use clustering methods to the same end include [Morinaga and Yamanishi \[2004\]](#) on tracking topic trends in real-time; [Stilo and Velardi \[2016\]](#) on temporal mining of microblog posts; and [Zhang et al. \[2015\]](#) on event detection and popularity prediction in microblogs.

Based on probabilistic topic models, [Blei and Lafferty \[2006\]](#) proposed a dynamic topic model to analyze time evolution of topics in document collections. In this model, documents are divided by time slice (year, month, etc) and, for each slice, a K-component topic model is obtained. As a result, the final topics are sequences of distributions over words, rather than a single distribution. With a similar purpose, [Wang and McCallum \[2006\]](#) presented Topics Over Time (TOT), a LDA-style model which uses both words co-occurrence and temporal information to discover topic changes across a period. Unlike the dynamic topic model, TOT avoids the discretization of time by calculating time stamp probabilities over all word tokens and obtaining topic distributions conditioned on a time stamp.

As a solution to the sparsity problem of topic detection on social networks, [Yin et al. \[2013\]](#) proposed a user-temporal mixture topic model to detect stable and punctual topics from social media data. Their framework uses the content of posts, their temporal information and the social network structure to identify and distinguish punctual topics from the stable ones. Their unified model significantly outperforms other approaches, which shows that using the social network and temporal structures can improve the process of temporal topic detection. Other works related to the discovering of temporal topics in social media, such as the studies of [Diao et al. \[2012\]](#) and [Xie et al. \[2016\]](#), focus on detecting and tracking bursty topics, i.e., they find topics that

cause a sudden increase in the number of posts within a short period of time and track their evolution.

In our study, we analyzed politicians tweets to find changes in their discussed topics over time. Following the idea of [Mei and Zhai \[2005\]](#) of comparing topics in consecutive time intervals by building an evolution graph, we proposed a method that uses BTM to discover topics on the short Twitter messages and builds a topic similarity graph to track variations on them across the studied period. Our model is explained in more details in [Chapter 4](#).

## 2.2 Polarity Evaluation

One of the key aspects to quantify polarization is to determine the polarity of each individual that belongs to the studied group. Polarity consists of the opinion of an individual regarding a topic, which she expresses by agreeing (“yes”), disagreeing (“no”) or being neutral towards it. In our study, we measure the polarity of a set of individuals as a previous step in the quantification of group polarization. For this reason, here we review some of the works regarding polarity evaluation.

The vast majority of the literature in polarity evaluation deals with the problem of classifying the position of an opinionated piece of text, which task is referred as *sentiment polarity classification* [[Pang and Lee, 2008](#)]. According to a detailed survey by [Ravi and Ravi \[2015\]](#), polarity classification has applications in many domains, such as the evaluation of product reviews, forums, blogs, news articles and micro-blogs. Identifying sentiments in micro-blogs such as Twitter, for example, is a big challenge due to the limits of characters and the informal language of the posts, which contains abbreviations and noisy words. Therefore, such texts need high level processing and more complex techniques for their analysis.

For the task of sentiment analysis in social media, some works use manually created sentiment resources or even create a set manually annotated words or posts for detecting sentiment in a dataset [[Wilson et al., 2005](#); [O’Connor et al., 2010](#); [Mohammad et al., 2017](#)]. Others explore some characteristics of the informal texts, such as hashtags and emoticons, to automatically find the polarity of the posts [[Davidov et al., 2010](#); [Kouloumpis et al., 2011](#); [Mohammad, 2012](#)]. The latter approach avoids the high cost of doing a manual annotation of the data and, for this reason, have been applied to analyze short informal texts [[Kiritchenko et al., 2014](#)]. As examples of studies that follow this approach, [Davidov et al. \[2010\]](#) and [Kouloumpis et al. \[2011\]](#) select hashtags that clearly show a sentiment – e.g., `#happy` or `#sad` – and use them as labels of



positive and negative sentiment to build a training dataset of tweets. The labeled data is used to train sentiment classifiers to find the polarity of tweets that do not contain the selected hashtags.

With the purpose of finding the polarity in a piece of text, the task of *stance detection* is also used by some recent studies to determine how favorable is an opinion regarding a target of interest [Mohammad et al., 2016]. Stance detection is related to some opinion mining techniques such as argument mining and sentiment classification, but it differs from the latter because it involves detecting the polarity towards a target that may not be explicitly mentioned in the text. There are a number of works of this task to perform target-dependent polarity classification in microblogs. Some of them use a manually annotated training dataset to find the stance of the tweets [Sobhani et al., 2016; Ebrahimi et al., 2016a; Taulé et al., 2017], others use a label propagation algorithm and a semi-supervised approach to classify the stance from the posts [Rajadesingan and Liu, 2014; Ebrahimi et al., 2016b].

In our study, one of our goals is to measure the polarity of the users, rather than solely identifying the opinion conveyed by the documents in a dataset. With the purpose of understanding the opinion dynamics of individuals in online social media, some works deal with the problem of modeling how users update their opinion in face of their neighbors' opinion and how their opinions evolve over time [Das et al., 2014; Morales et al., 2015; Jiang and Wu, 2017]. These studies determine the polarity of a set of users, which act as seeds of influence, and use opinion formation models to estimate the polarity of the other users in the network according to their interactions or relationships.

Also exploring the network structure, other works use collective classification to identify the opinions of users in social media [Li et al., 2016; Ileri and Karagoz, 2016]. Given a graph, where the nodes are users in a social network and the edges represent the relationships between them, and a subset of users which opinion is known (labeled users), Collective Classification deals with the problem of classifying the opinion of unlabeled users based on their links to the labeled ones. As another example of research that follows this direction, the work of Rabelo et al. [2012] investigate people opinion about the USA politics by building a directed graph, where the nodes are the social network users and the edges model the follower/follows relationship. As part of their method, they select a group of hashtags that clearly convey an opinion and classify posts according to the sentiment associated to these hashtags. In order to find a set of users with a known opinion, users are labeled with the opinion that has the highest count among their posts. The collective classification is then performed on users who have no posts. The authors recorded a precision of 80% when they have as few as 10%

of labeled users in the network, which means that their approach is effective to classify at least 90% of the users in the graph.

Following a similar idea to [Rabelo et al. \[2012\]](#), our study finds the polarity of users in social media by building a retweet network graph, where a relationship between two users exists when one retweets the other. In our approach, we selected hashtags that clearly indicate a political lean and classified the polarity of posts according to these hashtags. We then label a sample of users according to their posts, and subsequently examine the relationships in the retweet network to find the polarity of the unlabeled users. Our methods to evaluate polarity of users in social media were used to study of polarity of the Brazilian general public (Chapter 5).

## 2.3 Polarization Measurement

Polarization consists of a social phenomenon characterized by the simultaneous presence of people that hold divergent opinions or principles. Studying polarization is relevant to understand conflicts in society and the evolution of differences in the beliefs of sets of individuals. Besides, investigating this phenomenon makes it possible to predict and even minimize tensions among social groups. For this reason, there are a number of works that deal with the problem of measuring polarization, which is our main goal of this research. We describe some of these studies in this section.

Although there are many works in the field of opinion polarization, there is no consensus about a quantitative measure for it [[Schmitt, 2016](#)]. According to [Bramson et al. \[2016\]](#), most of the studies present a formal measure for polarization which is specific to the dataset or topic of interest (e.g., politics), which explains the diversity of polarization measures in the literature.

With the purpose to quantify controversy, [Garimella et al. \[2018\]](#) propose the Random Walk Controversy (RWC). Given two partitions  $A$  and  $B$ , this measure basically quantifies the likelihood of an individual, which is in a partition  $A$ , to be exposed to authoritative content from the opposite partition (partition  $B$ ). RWC is independent from the size of partitions, as well as the degree of the vertices in each of the partitions. A high RWC indicates a low probability of the individual to be exposed to the contrary content, demonstrating that the opposite sides are highly closed in their circles. On the other hand, a low RWC indicates that the probability of being exposed to the opposite content is similar to the probability of finding content from the same side, which show that the individuals from the different partitions are more open to the contrary points of view.

Many works, such as [Conover et al. \[2011\]](#) and [Guerra et al. \[2013\]](#), use the modularity of a graph to measure the polarization of individuals. Proposed by [Newman \[2006\]](#), modularity quantifies the level of division of the network into groups, i.e., the extent to which the nodes of a graph are grouped into different communities. Therefore, a high modularity value indicates that the network has dense connections between the nodes within the groups, but few connections between nodes from different groups.

[Morales et al. \[2015\]](#) present the polarization index, a measure which takes into account the probability density distribution of the opinions of individuals to quantify the segregation within a population. The measure considers that group polarization depends on the difference between the size of the opposite groups, as well as the distance between their central points of view. Even though the authors use a network to estimate the polarity of individuals in their study, they do not need a network structure to calculate the final polarization, which is solely based on the density distribution of opinions.

We used the polarization index proposed by [Morales et al. \[2015\]](#) to quantify polarization in our study. Besides the solid social fundamentals behind it, we chose this measure due to the fact that it does not require a network structure to compute polarization, which overcomes the restrictions to calculate polarization in our real-world dataset. In view of its crucial relevance in our study, [Chapter 3](#) gives a more detailed explanation of the polarization index.

## 2.4 Computational Social Science and Politics

In view of its growing popularity among people around the globe, online social networks are often mirrors of what society thinks. For this reason, they have a strong potential in the study of social preferences and actions, as well as in predicting events and avoiding problems in the real world. In virtue of these possibilities, recent research explore social networks to understand the links between virtual and real-world actions in a variety of applications, such as the study of posts about mental health in online communities to predict suicide ideation [[De Choudhury et al., 2016](#)] or the analysis of students' informal conversations on social media to understand problems in education systems [[Chen et al., 2014](#)]. These works are part of the literature on Computational Social Science, an interdisciplinary field that studies the dynamics of society with the aid of computational methods [[Cioffi-Revilla, 2013](#)]. It involves studies that collect and process data so as to investigate patterns of individual and group behaviors, aiming to understand important aspects of society. According to [Lazer et al. \[2009\]](#), social media

platforms can reveal a complete record of individual behavior, offering opportunities to understand the impact of a person’s position in the network and the changes to this position over time.

One of the most studied applications of Computational Social Science is politics, since the Web and the social platforms are valuable tools to analyze the views and beliefs of a population and predict outcomes regarding political events. As an example of these studies, [Morales et al. \[2015\]](#) investigated the emergence of polarization towards the ex-Venezuelan president Hugo Chávez on Twitter. By analyzing a set of retweets, their studies showed that a small set of influential users was able to spread their opinions through social networks and generate an impact in the segregation of opinions of the population. Moreover, they compared the social media data to socio-economic offline data, which results showed that online polarization is closely related to political, geographical and social polarization in Venezuela.

[Conover et al. \[2011\]](#) also studied political polarization on Twitter by building retweet and mention networks from a dataset of tweets regarding the 2010 presidential elections in the United States. They used clustering algorithms to process these networks so as to understand people’s interactions concerning the elections on social media. Their results showed that retweet networks have a more segregated structure, since users retweet each other to endorse an opinion. On the other hand, mention networks have a more diverse set of users and are characterized by the interaction of users of opposite political views. Based on this same purpose, some studies also measure divergence between political groups by detecting and analyzing communities in social networks [[Guerra et al., 2013](#); [Adamic and Glance, 2005](#)].

Also exploring the structure of social networks, [Cota et al. \[2019\]](#) analyzed the polarization around the impeachment of the Brazilian ex-president Dilma Rousseff on Twitter. Their study involved building a mention network to capture real social interactions regarding the event. To calculate polarization, they also classified the tweets exchanged among the users by manually labeling a group of related hashtags, which allowed them to quantify the “echo chambers” in the network.

Apart from evaluating the opinions from the general public, other studies focus on the interactions of political elites on social media. [Lietz et al. \[2014\]](#), for instance, proposed a set of quantitative measures to study the socio-cultural structure and dynamics of the online conversational practices of political parties on Twitter over time. Following the same line, [Livne et al. \[2011\]](#) used graph and text mining techniques to study the behavior of the candidates of the most popular parties in the United States during the 2010 American elections. By analyzing their tweets, the authors identified noticeable differences on how each candidate interact on social media. Besides, the

textual aspects of their posts and the graph structures for each party were used to create a model to predict the outcomes of the elections, which recorded an accuracy of 88%.

Our work builds on this body of research by proposing computational methods to analyze online and offline data in order to study political polarization. We investigate real-world and social media data from politicians and the general public regarding the Brazilian political scenario and the impeachment proceedings of Dilma Rousseff in 2016, for which we also analyze possible associations between online aspects and offline actions. The following chapters describe the methods and the experimental results of our research.

# Chapter 3

## Polarization Index

In Social Sciences, polarization consists of a social phenomenon that raises contrary reactions from a set of individuals, causing them to move towards an extreme point according to their own previous tendencies [Sunstein, 2002; Dixit and Weibull, 2007]. Its concept is related to a group behavior, which is characterized by the simultaneous presence of people that hold conflicting opinions or principles, i.e., the co-existence of sets of individuals having opposite polarities. The polarity of a person, in turn, is defined here as her position regarding a subject, which she expresses by agreeing (“yes”), disagreeing (“no”) or being neutral when confronted with the topic.

As explained in Section 2.3, there are many ways to measure polarization. We chose to use the polarization index proposed by Morales et al. [2015], because this measure does not require a network structure to compute polarization, which overcomes the restrictions to calculate polarization in our real-world dataset. Besides, the measure has solid social fundamentals behind it, since it covers the 3 axioms of polarization [Esteban and Ray, 1994]: I) a polarized population has a high degree of homogeneity within each group; II) a polarized population has a high degree of heterogeneity across groups; and III) in the population, there must be a small number of groups of different opinions with a significant size.

Since understanding the polarization index is crucial to describe the methods and results of our research, the following sections present the variables related to the final polarization metric: the difference between populations sizes (Section 3.2), the distance between gravity centers (Section 3.3) and, finally, the polarization index (Section 3.4).

### 3.1 Background

First of all, [Morales et al. \[2015\]](#) consider that a population is perfectly polarized when it is divided into two groups of the same size that share opposite views about a subject. This definition is based on the idea of the electric dipole moment, which measures the overall polarity of a charge system. A dipole is created by the combination of two opposite charges that have equal magnitude and are separated by a certain distance. Hence, considering the basic case of a dipole, the electric dipole moment is proportional to the distance between the charges, that is, it increases with the separation of the charges. The authors borrowed this notion to describe the phenomenon of opinion polarization: as happens to the dipole, the polarization of two different groups depends on the distance between their points of view, i.e., how different their opinions are.

Figure 3.1 illustrates variables involved in the computation of the polarization index. The metric takes into account the size of the populations of opposite opinions ( $A^-$  and  $A^+$ ), the gravity centers of each population ( $gc^-$  and  $gc^+$ ) and the distance between these gravity centers ( $d$ ). The following sections describes the mathematical definition for these variables, introducing their formulas and the reasoning behind them.

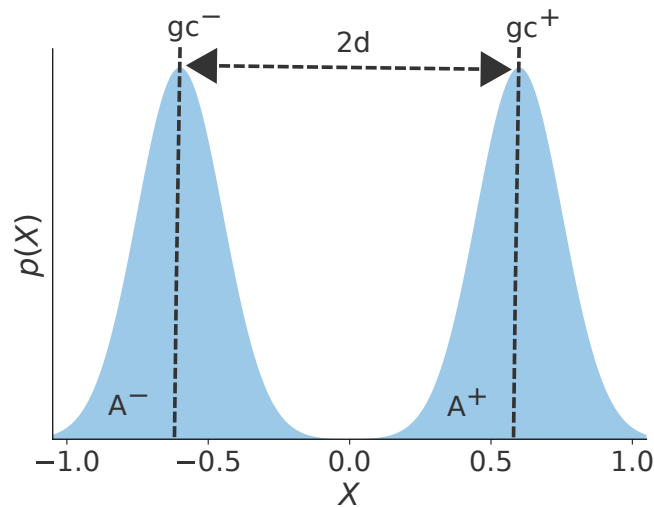


Figure 3.1: Example of probability density distribution of polarities, showing the populations of opposite opinions ( $A^-$  and  $A^+$ ), the gravity centers of each population ( $gc^-$  and  $gc^+$ ) and the distance between these gravity centers ( $d$ ).

Adapted from Morales [Morales et al. \[2015\]](#).

## 3.2 Difference between populations sizes

Compared to the charges of a dipole, the groups or populations are here regarded as set of individuals that hold an opinion, having a polarity value  $X$  associated to each member. Given that polarity  $X$  is measured for each individual in the range  $[-1, 1]$ , it is possible to calculate the size of such populations by taking into account the probability density distribution of polarities  $p(X)$  for the set of studied individuals. Thus, the population of negative opinions ( $X < 0$ ), represented by  $A^-$ , is computed by integrating the distribution  $p(X)$  over the interval  $[-1, 0]$  (Equation 3.1). Comparatively, the population of positive opinions ( $X > 0$ ), designated by  $A^+$ , is calculated by integrating the polarity distribution over the interval  $[0, 1]$ , as shown in Equation 3.2.

$$A^- = \int_{-1}^0 p(X) dX = P(X < 0) \quad (3.1)$$

$$A^+ = \int_0^1 p(X) dX = P(X > 0) \quad (3.2)$$

Having calculated these values, Equation 3.3 determines the normalized difference between population sizes  $\Delta A$ , which is one of the central variables to compute the final polarization index. This difference represents how unbalanced are the existing groups, i.e., it shows if one population has a greater density of individuals than the other.

$$\Delta A = |A^+ - A^-| \quad (3.3)$$

Note that the population sizes  $A^-$  and  $A^+$  are calculated as the area under the probability density function over the polarities which, respectively, represent negative and positive opinions. Therefore, their values reveal the probability of an individual to be part of that population. As probabilities,  $A^-$  and  $A^+$  values lie in the range  $[0, 1]$  and, as a result, their normalized difference  $\Delta A$  also is restricted to the range between 0 and 1. The closer  $\Delta A$  is to 0, the more similar the population sizes. Conversely, a  $\Delta A$  close to 1 indicates that the probability distribution takes the shape of an unimodal distribution, having one population with a much greater density than its opposite one.

## 3.3 Distance between gravity centers

Another key variable is the distance  $d$  between the positive and negative opinions, which quantifies the level of divergence between the opposite populations. It takes into account the gravity centers of negative  $gc^-$  (Equation 3.4) and positive opinions



$gc^+$  (Equation 3.5), which measure the central opinion of the positive and negative populations.

$$gc^- = \frac{\int_{-1}^0 p(X)X dX}{\int_{-1}^0 p(X) dX} \quad (3.4)$$

$$gc^+ = \frac{\int_0^1 p(X)X dX}{\int_0^1 p(X) dX} \quad (3.5)$$

The distance  $d$  is then computed as the normalized difference between these gravity centers, as shown by Equation 3.6:

$$d = \frac{|gc^+ - gc^-|}{|X_{\max} - X_{\min}|} = \frac{|gc^+ - gc^-|}{2} \quad (3.6)$$

In this Equation,  $X_{\max}$  represents the upper limit of the opinion values of the positive population (i.e.,  $X_{\max} = 1$ ) and  $X_{\min}$  represents the lower limit of the opinion values of the negative population (i.e.,  $X_{\min} = -1$ ).

Note that  $d = 0$  indicates that the individuals share the same opinion, since there are no difference between the central opinions of the opposite populations. On the other hand, a distance  $d$  close to 1 reveals that the two main opinions are in the extremes of each side.

### 3.4 Polarization Index

After computing the previous variables, Equation 3.7 shows how to compute the polarization index  $\mu$ , which is finally calculated as a function of the difference between populations sizes  $\Delta A$  and the distance between the gravity centers  $d$ . As previously explained, polarization increases with the separation of the opposite groups, which is why the index  $\mu$  is proportional to the distance between gravity centers  $d$ . Also, polarization is affected by the density of the populations, reaching its maximum value when the groups have equal sizes. Alternatively, the greater the difference between these groups, the smaller the index  $\mu$ .

$$\mu = (1 - \Delta A) d \quad (3.7)$$

The polarization index  $\mu$  lies in the range  $[0, 1]$  and its resulting values can be interpreted as follows. When  $\mu$  reaches its maximum value ( $\mu = 1$ ), we can say that the population is perfectly polarized. In this case, the populations have equal sizes

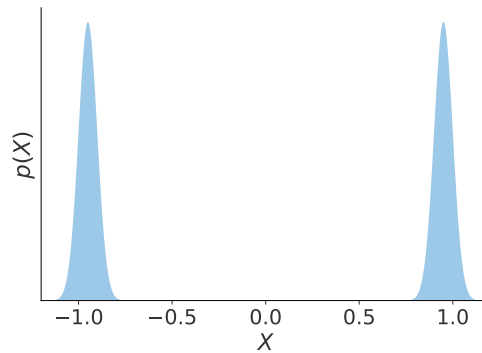


Figure 3.2: Example of probability density function having maximum polarization index ( $\mu = 1$ ).

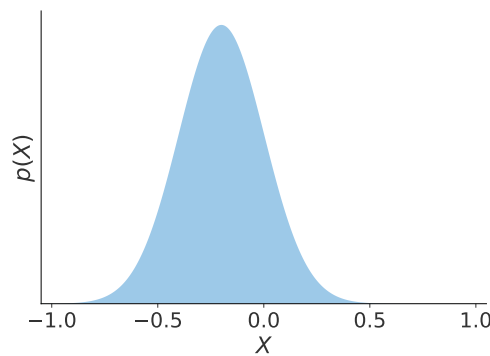


Figure 3.3: Example of probability density function having a population that is not polarized ( $\mu = 0$ ).

and their polarities are centered in the extreme values ( $-1$  and  $1$ ), as exemplified by Figure 3.2.

When  $\mu$  has its minimum value ( $\mu = 0$ ), it means that the population is not polarized. This is the case that the probability distribution of polarities takes the shape of an unimodal distribution, having the difference between populations sizes  $\Delta A = 1$ . In this situation, either the population is centered at a neutral opinion or it is entirely centered in one of the extremes. Figure 3.3 shows an example of this case.

The polarization index falls in the limits of its interval due to a combination of variations in its related variables  $\Delta A$  and  $d$ . Morales et al. [2015] describe three main situations which may lead to it:

- The populations have the same size ( $\Delta A = 0$ ), but the distance between their gravity centers  $d$  is smaller than 1. In this situation, the whole group of people is equally divided into two sets of individuals, each of which supports an opinion

that is opposite to the other. The gravity centers of these opinions, however, are less opposed to each other when compared to the perfectly polarized scenario. As a result, the level of divergence between the central opinions – represented by the distance between the gravity centers  $d$  – dictates the final polarization index  $\mu$ . The initial Figure 3.1 shows an example of this case.

In this scenario, the polarization index does not reach its maximum ( $\mu < 1$ ) due to the reduced distance between the gravity centers. It demonstrates that, even if a set of individuals is equally divided by its conflicting points of view, group polarization depends on how contrasting are the central opinions of each of the two opposed groups.

- The distance between gravity centers  $d$  has its maximum value ( $d = 1$ ), but the population sizes are different ( $\Delta A > 0$ ). In this case, one of the opposite populations attracts a greater density of individuals than the other one, i.e., one of the populations has a majority of supporters of its opinion. The gravity centers of each population are located in the extremes ( $gc^- = -1$  and  $gc^+ = 1$ ), which means the distance between them achieves its maximum value ( $d = 1$ ). Therefore, the polarization index  $\mu$  is given by the difference between the populations sizes  $\Delta A$ . Figure 3.4 shows an example of this situation.

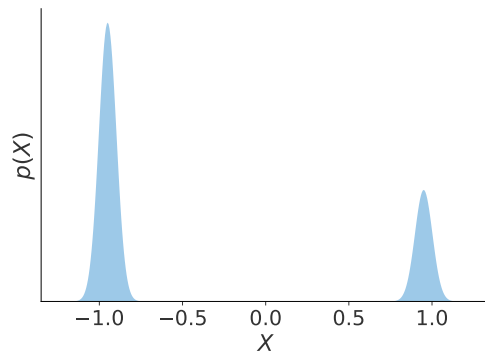


Figure 3.4: Example of probability density function where  $d = 1$  and  $\Delta A > 0$ .

Although the central opinions are totally opposed in this scenario, the polarization index does not reach its peak ( $\mu < 1$ ). In this way, the metric indicates that, the larger the population that concentrates one of the opinions, the smaller the group polarization. To put it differently, if the most of the individuals hold a certain opinion, it tends to lessen the level of polarization of the group, even though there are some people in the population that are completely contrary to the position of the majority.

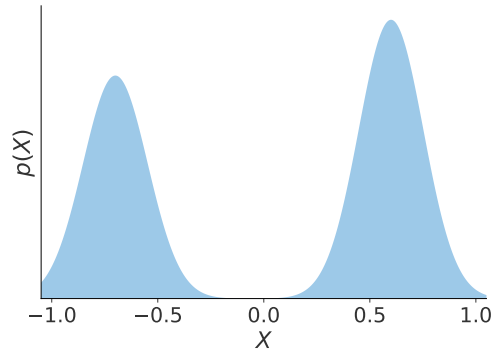


Figure 3.5: Example of probability density function where  $0 < d < 1$  and  $0 < \Delta A < 1$ .

- Both the difference between populations sizes  $\Delta A$  and the distance between gravity centers  $d$  fall in the limits of their ranges. This is the situation where one of the populations have a greater density of individuals than the other, and the gravity centers of each population are not located in the extremes ( $-1 < gc^- < 0$  and  $0 < gc^+ < 1$ ). One of the examples that illustrate this situation is shown on Figure 3.5.

In this case, polarization index  $\mu$  is a combination of the value of  $\Delta A$  and  $d$ . Hence, in order to identify which factor is the most relevant to determine group polarization, one can investigate the temporal evolution of the index  $\mu$  and its variables. By comparing the values of polarization to its related variables over time, it would be possible to understand which of the variables has more influence over the changes in the polarization index.

In short, the chosen metric takes into account the difference between the central opinions – represented by the distance between the gravity centers – and the size of the divergent groups – described by the difference between the population sizes. For this, it reflects the basic aspects of a polarized set of individuals in the literature: a small number of highly segregated groups of similar size, which ideas inside a group are extremely close. As the polarization axioms postulate, the more homogeneous are the ideas of the individuals inside a group and the more heterogeneous are the opinions of the groups, the larger the polarization. Besides, the closer the sizes of the opposed groups, the larger the population.

# Chapter 4

## Polarization of Politicians

In the past few years, the revelation of corruption scandals, such as the *Car Wash* operation, made public a long list of Brazilian politicians which were accused of being involved in a complex bribery scheme in Petrobras, the Brazilian largest state-oil company. The investigation included politicians from the many different parties in the Brazilian political scenario, some of which were led to prison. Among the investigated politicians, there were members of the Workers Party (PT), party to which belongs Dilma Rousseff, the president in charge at that time. Besides causing tension and disbelief in the citizens, this scenario also aroused conflicts between pro and anti-government parties, intensifying the political instability in the Congress. In view of these circumstances, opposition parties and movements organized protests against the government [Venceslau, 2016]. These recurrent turbulences in and out of the Congress culminated in the impeachment of Rousseff in 2016.

Considering this context, this Chapter presents an study of the polarization phenomenon among Brazilian politicians. In our study, we investigate the behavior of these politicians in social media by evaluating how their main discussed topics change over time (Section 4.1). In addition, we quantify the polarization among them by analysing their roll-call votes in bills at the Lower House (Section 4.2). At the end, we compare their most relevant topics to their polarization measures over time, in order to understand possible associations between what they say in social media and their voting behavior in the Congress (Section 4.3).

### 4.1 Investigating Topic Evolution in Social Media

In this work, one of our main goals is to study the behavior of Brazilian politicians on the Web by applying algorithms to analyze and summarize information from their

social media data. To achieve this, we collected a dataset of Twitter posts from Brazilian representatives (Section 4.1.1) and proposed a method to investigate how their discussed topics change over time (Section 4.1.2). Our results showed that the political crisis and the activities at the Lower House are the mostly discussed topics in social media among the representatives. The political crisis topic covers almost the entire period, being even more intensely discussed during the months associated with the impeachment proceedings (Section 4.1.3).

### 4.1.1 Data

We collected tweets from Brazilian representatives that are part of the Lower House of Congress (House of Representatives) and have an active Twitter account. The dataset includes 502,342 tweets from 423 representatives (about 82.5% of the total numbers of the House) shared between January of 2015 and November of 2016. As a preprocessing step, all tweets were lower-cased and stop words were eliminated. All messages in our Twitter dataset are in Portuguese, but our results were translated to English for the sake of understanding.

### 4.1.2 Methods

Although there are a few works in the literature for temporal topic modeling (Section 2.1), our study is concerned with detecting both stable and temporal topics, rather than discovering bursty patterns, which is the goal of many works regarding temporal topic detection on social media [Diao et al., 2012; Xie et al., 2016]. Apart from that, our politicians data is able to generate a quite small user network structure, comprising only 15% of the dataset, which imposes a restriction on using user-topic mixture models, such as the one proposed by Yin et al. [2013]. Thus, here we proposed a method that follows the basic idea of Mei and Zhai [2005] of comparing topics in consecutive time intervals by building an evolution graph. Unlike Mei and Zhai [2005], our approach compares topics from all the time slices and measures the similarity between them. Besides, we take into consideration that sets of similar topics may be related to others, for which we successively group the sets according to a similarity threshold. In order to deal with the restrictions of our study, our methodology uses the Biterm Topic Model (BTM) to discover topics from fixed time slices and builds a topic similarity graph to track variations on them across the studied period.

As explained in the Section 2.1, the Biterm Topic Model (BTM) finds topics by modeling the generation of biterms in a corpus. Figure 4.1 shows the graphical

representation for the algorithm. BTM considers the whole collection of documents as a single one, modeling the corpus as a mixture of topics. This algorithm draws a topic distribution  $\theta$  for the entire collection and, for each topic  $z$ , it draws a topic-specific word distribution  $\phi$ , according to the probability of the co-occurring words in the corpus.

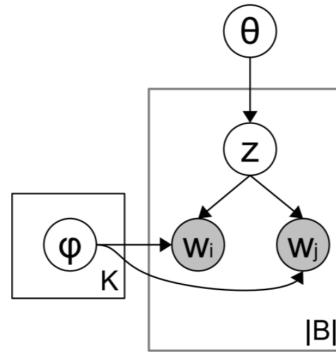


Figure 4.1: Graphical Representation of Biterm Topic Model (BTM).  
Extracted from Yan et al. [2013].

Our temporal topic evolution method works as follows. To compose the dataset of interest, a set of messages  $M$  is extracted from Twitter over a period of time  $P$ . These messages are organized according to a time slice of interest (e.g., hour, day, month, etc)  $M_P = \{M_1, M_2, \dots, M_p\}$ , where  $p$  is the number of time slices in  $P$ . In our study, we break the studied period into monthly slices. Then, a topic modeling algorithm is used to extract topics from every set  $M_i$ . Here we work with BTM (Biterm Topic Modeling), since it is considered a state-of-art method for extracting topics from short texts [Cheng et al., 2014]. It receives as a parameter the number  $k$  of topics and a time slice  $i$ , and produces as output a list of topics  $T_i = \{T_i^1, T_i^2, \dots, T_i^k\}$ , where each  $T_i^j$  is represented by its topic number  $j$  (i.e., the number that identifies a topic) and is described by  $w$  words. After executing BTM for every time slice, we end up with a set of topics  $T_P = \{T_1, T_2, \dots, T_p\}$  for the entire period. This Topic Modeling step is described by the Algorithm 1.

Figure 4.2 illustrates the core steps of the proposed method. From  $T_P$ , we create a unique topic similarity graph  $G_T = \{V_T, E_T\}$  to find groups of similar topics, as described by the Algorithm 2. In this graph, the vertices  $V_T$  represents the topics in  $T_P$  and the edges in  $E_T$  measure the similarity between each pair of topics. The similarity was calculated using the Jaccard coefficient  $\sigma$ , which computes the proportion of shared words  $w$  between each pair of topics  $T_{i_1}^{j_1}$  and  $T_{i_2}^{j_2}$ :

**Algorithm 1** Temporal Topic Evolution Algorithm - Part 1 - Topic Modeling**Data** $M$  set of Twitter messages divided according to the chosen time slice**Input** $p$  number of time slices $k$  number of topics**Output** $T_P$  set of topics for each time slice $T_P \leftarrow \emptyset$ 

▷ Initialize an empty set of topics

**for**  $i \leftarrow 1, p$  **do** $T_i \leftarrow \text{BTM}(M_i, k)$ 

▷ Use BTM to find the topics for each time slice

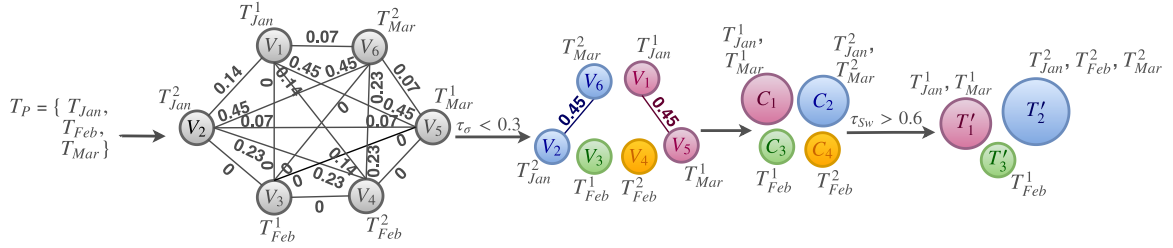
INSERT( $T_i, T_P$ )**end for**

Figure 4.2: Example of application of the proposed topic evolution method.

$$\sigma(T_{i_1}^{j_1}, T_{i_2}^{j_2}) = \frac{|w_{T_{i_1}^{j_1}} \cap w_{T_{i_2}^{j_2}}|}{|w_{T_{i_1}^{j_1}} \cup w_{T_{i_2}^{j_2}}|} \quad (4.1)$$

In the example illustrated in the figure, we consider three monthly slices, for which we have the sets of topics  $T_P = \{T_{Jan}, T_{Feb}, T_{Mar}\}$ . These sets of topics have two topics each, resulting in a graph with 6 nodes. Note that each node corresponds to a topic  $T_i^j$ , where  $j$  is the topic number and  $i$  is its associated monthly slice.

From  $G_T$ , we want to find groups of similar topics and, to achieve that, two other steps are followed. First, we remove from  $E_T$  the edges that indicate a similarity lower than a threshold  $\tau_\sigma$  (in Figure 4.2, for example,  $\tau_\sigma = 0.3$ ). After that, we end up with a set of  $m$  connected components  $C_T = \{C_1, C_2, \dots, C_m\}$ , which we assume to capture an intrinsic similarity between the topics related to the event of interest. In this way, each  $C_x$  is considered as a group of similar topics, and the number of topics is reduced from the original  $k \times p$  to  $m$ . We merge all the nodes (topics) within  $C_x$  into a single super-node and, consequently, a super-topic, which is now described by the union of the words of each topic belonging to  $C_x$ . Note that  $C_x$  may have components from



---

**Algorithm 2** Temporal Topic Evolution Algorithm - Part 2 - Topic Similarity Graph

---

**Input** $T_P$  set of topics for each time slice**Output** $G_T$  Topic Similarity Graph $G_T \leftarrow \{V_T, E_T\}$ ▷ Initialize the Topic Similarity Graph  $G_T$  $V_T \leftarrow T_P$ ▷ Each topic is a vertice in  $G_T$  $E_T \leftarrow \emptyset$ **for all** pair of vertices  $(v_1, v_2) \in V_T$  **do** $e_\sigma \leftarrow \frac{|w_{v_1} \cap w_{v_2}|}{|w_{v_1} \cup w_{v_2}|}$ 

▷ Jaccard's Similarity

Add edge  $(v_1, v_2)$  with weight  $e_\sigma$  to  $E_T$ **end for**Remove edges which  $e_\sigma < \tau_\sigma$  from  $E_T$  ▷ No connections between dissimilar topics

---

different time slices, and that is our goal: to say that the super-topic representing  $C_x$  appears in different time slices. For instance, Figure 4.2 shows that  $C_1$  contains the topic number 1 of January ( $T_{\text{Jan}}^1$ ) and the topic number 1 from March ( $T_{\text{Mar}}^1$ ), meaning that this super-topic is discussed in January and March.

Our original graph  $G_T$  has now a set of super-nodes  $N_G = \{C_1, C_2, \dots, C_m\}$ . However, after each connected component becomes a topic, these new topics may again share a large number of words  $w$  with others, and hence could be merged again. As a result, we refine these super-nodes representing topics by successively merging them. This is done by calculating, for each pair of nodes  $(C_i, C_j)$ , their percentage of shared words Sw:

$$\text{Sw}(C_i, C_j) = \frac{|w_{C_i} \cap w_{C_j}|}{\min(|w_{C_i}|, |w_{C_j}|)} \quad (4.2)$$

The condition to merge the vertices is the following: if Sw is larger than a threshold  $\tau_{\text{Sw}}$  (in Figure 4.2, e.g.,  $\tau_{\text{Sw}} = 0.6$ ), the pair of vertices  $(C_i, C_j)$  is grouped. This grouping process continues until there is no pair of vertices that meets the condition. Algorithm 3 shows this successive grouping step.

At the end of the grouping process, we have a smaller number of  $n$  super-topics  $T'_P = \{T'_1, T'_2, \dots, T'_n\}$ , each one with its associated time slices, which allow us to follow the evolution of these topics over time. Since each  $T'_x$  is composed by a group of connected components, it is also represented by the original topics that are included in these components and their words. As shown in Figure 4.2, the initial 6 original topics

---

**Algorithm 3** Temporal Topic Evolution Algorithm - Part 3 - Successive Grouping

---

**Input** $G_T$  Topic Similarity Graph**Output** $T'_P$  set of final topics**do** $V'_T \leftarrow \emptyset$ **for all** connected component  $c \in V_T$  **do** $v' \leftarrow \emptyset$ **for all** vertice  $v \in c$  **do**INSERT( $w_v, v'$ ) $\triangleright$  The words of  $c$  are merged into vertice  $v'$ **end for**INSERT( $v', V'_T$ ) $\triangleright$  New vertice  $v'$  is added to the set  $V'_T$ **end for** $V_T \leftarrow V'_T$  $E_T \leftarrow \emptyset$ **for all** pair of vertices  $(v_1, v_2) \in V_T$  **do** $e_{Sw} \leftarrow \frac{|w_{v_1} \cap w_{v_2}|}{\min(|w_{v_1}|, |w_{v_2}|)}$  $\triangleright$  Successive Grouping Similarity**if**  $e_{Sw} > \tau_{Sw}$  **then**Add edge  $(v_1, v_2)$  with weight  $e_{Sw}$  to  $E_T$ **end for****while**  $G_T$  is a connected graph $T'_P \leftarrow V_T$  $\triangleright$  The final topics are the connected components of the graph  $G_T$ 

---

$(T_{Jan}^1, T_{Jan}^2, T_{Feb}^1, T_{Feb}^2, T_{Mar}^1, T_{Mar}^2)$  were grouped into 3 final super-topics  $(T'_1, T'_2, T'_3)$ . In this example,  $T'_2$  contains the topics  $T_{Jan}^2, T_{Feb}^2$  and  $T_{Mar}^2$ , meaning that this super-topic is discussed in January, February and March.

Following the generation of the final super-topics  $T'_P$ , we want to quantify their relevance as well as the relevance of the words describing them. These metrics of relevance are based on word and topic probabilities which are calculated by the topic modeling method (BTM), since BTM produces as output the probability of each topic and the probability of a word given a topic, as explained in Section 2.1. We also use the BTM algorithm to find the proportion of Twitter messages for each original topic  $T_i^j$ , since this algorithm already assigns the most probable topic to each message in the dataset.

The topic relevance  $TR_x^i$  measures the popularity of a final super-topic  $T'_x$  in a time slice  $i$ . It is calculated as the sum of the total of messages  $M$  assigned to each

original topic  $T_i^j$  that belongs to the super-topic  $T'_x$ , as shown in the following Equation:

$$\text{TR}_x^i = \sum_{T_i^j \in T'_x} |M_{T_i^j}| \quad (4.3)$$

We also compute word relevance for each word  $w$  in the vocabulary and each super-topic  $T'_x$  at a time slice  $i$ . It is important to remember that each super-topic is composed by a set of original BTM topics  $T_i^j$  at different times. Hence, here we use  $T'_{x,i}$  to represent the super-topic  $T'_x$  at time slice  $i$ , which only comprises the set of original topics  $T_i^j$  found at time  $i$ . Having that in mind, word relevance  $\text{WR}_{x,i}^w$  is given by  $P(w | T'_{x,i})$ , i.e., is calculated as the probability of a word  $w$  given a super-topic  $T'_x$  at a time slice  $i$ , as:

$$\text{WR}_{x,i}^w = P(w | T'_{x,i}) = \sum_{T_i^j \in T'_{x,i}} P(w | T_i^j) \cdot P(T_i^j) \quad (4.4)$$

The formula for word relevance is based on the probability of the original topic  $P(T_i^j)$  and the conditional probability of a word given the original topic  $P(w | T_i^j)$ , both previously computed by the BTM.

### 4.1.3 Experimental Results

In order to understand what the politicians were discussing on Twitter, we first characterized the dataset in terms of the number of tweets over the months, as shown in Figure 4.3. As can be observed, representatives significantly increased their participation on Twitter in 2016: there was no month that registered more than 20K posts in 2015, whereas all months exceeded this number after February 2016.

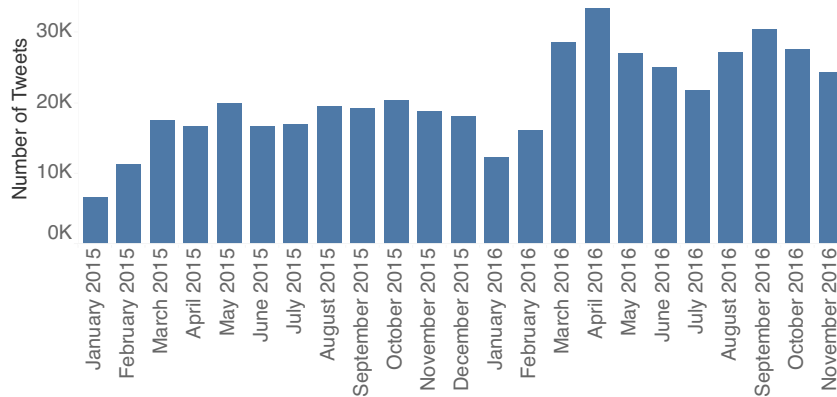


Figure 4.3: Number of tweets over time for the politicians.

In addition, we also investigated how the participation of the politicians on Twitter changed over time, according to the political spectrum of the parties to which they are part of (Figure 4.4). In this analysis, the participation  $P_s$  is given by the total of tweets  $M$  posted by the politicians of some political spectrum  $s$  divided by the number of the representatives  $R$  which are members of parties from that spectrum:

$$P_s = \frac{|M_s|}{|R_s|} \quad (4.5)$$

As shown in Figure 4.4, the politicians of the left-wing spectrum participated more on Twitter than politicians of the other spectra for most of the studied period, having a major increase in their participation from March of 2016 onwards. We can also observe that politicians of the other political spectra, except for the right-wing, increased their participation in 2016, although in a smaller proportion than the left-wing ones. It is also important to notice that the politicians of all spectra have a peak in their participation on Twitter on April 2016, when the impeachment voting took place at the Lower House. More especially, the larger increase in the participation of the politicians of the left-wing spectrum can be associated with that impeachment voting event, since the ex-president Dilma Rousseff is a member of a left-wing party.

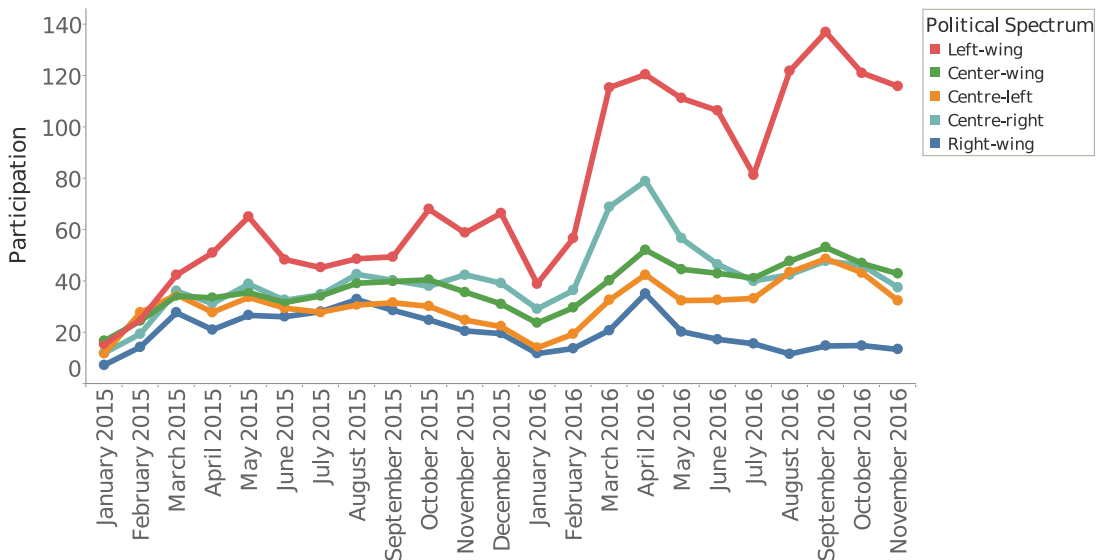


Figure 4.4: Participation of the politicians on Twitter according to the political spectrum of their parties ( $P_s$ ) over time.

After data characterization, we followed the proposed method for temporal topic evolution (Section 4.1.2) to investigate which topics were discussed over the period and how long the topics last. In this analysis, we considered a month as a time unit. We set the value of the parameter  $k$  of BTM as 10 and, over the 23 months of analysis

( $p = 23$ ), we obtained a total of 230 topics. Each of these topics are defined by the 10 most probable words returned by BTM. It is also important to say that the thresholds  $\tau_\sigma$  and  $\tau_{Sw}$  were defined by a qualitative analysis of the intermediary results, due to restrictions of measuring the topic coherence for Portuguese language. The qualitative analysis showed that the thresholds  $\tau_\sigma = 0.35$  and  $\tau_{Sw} = 0.80$  yield the most coherent results for our method.

Figure 4.5 shows word clouds for the original topic models from four months of our dataset: April 2015, August 2015, April 2016 and August 2016. Notice that the top words differ among months: in April and August 2015, “House” (*House of Representatives*), “Today”, “Day” and “Representative” are the most prominent words; but “Impeachment”, “Dilma” (*Dilma Rousseff*) and “Brazil” are the most frequent ones in the same months of the following year. This difference can be explained by the fact that impeachment proceedings were officially launched on December 2015 at the House of Representatives (see Section 1.1). Although there were protests from people to demand the impeachment of Dilma over 2015, the process was only discussed and voted at the Lower House on April of 2016. For this reason, the politicians seem to have given more emphasis to the impeachment subject during the year of 2016, whereas they discussed apparently common matters before December of 2015.

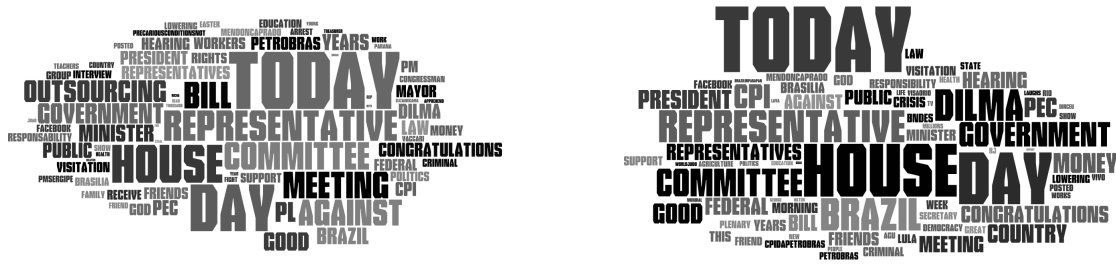
Following our method, we built the topic similarity graph and used it to find a set of super topics. From the 230 initial topics (10 topics for each of the 23 months), 50 super topics were obtained after the aggregation process ( $n = 50$ ), and their relevance was calculated per month. Results are summarized in the heat map of Figure 4.6, which is filtered to show the 5 most relevant aggregated topics for simplification<sup>1</sup>. Remember that topic relevance  $TR_i^x$  is computed by the number of tweets posted in a month, and its value is color-coded using a log scale in the heat map. To understand the content of these topics, the top 10 words for each final topic are also presented in Table 4.1.

According to the top words of each final topic, we can make the following observations:

- $T'_1$  contains words that are closely related to Brazilian political crisis (“impeachment”, “coup”, “against”).
- $T'_2$  seems to comprise posts about the activities of the politicians at the Lower House, since it includes words such as “committee”, “meeting” and “bill”.

---

<sup>1</sup>The complete interactive heat map can be seen on <http://homepages.dcc.ufmg.br/%7Eroberta.coeli/mestrado/TopicsRelevanceOverTime.html>



(a) April 2015 ( $T_{Apr/2015}$ )

(b) August 2015 ( $T_{Aug/2015}$ )



(c) April 2016 ( $T_{Apr/2016}$ )



(d) August 2016 ( $T_{Aug/2016}$ )

Figure 4.5: Word clouds for monthly datasets.

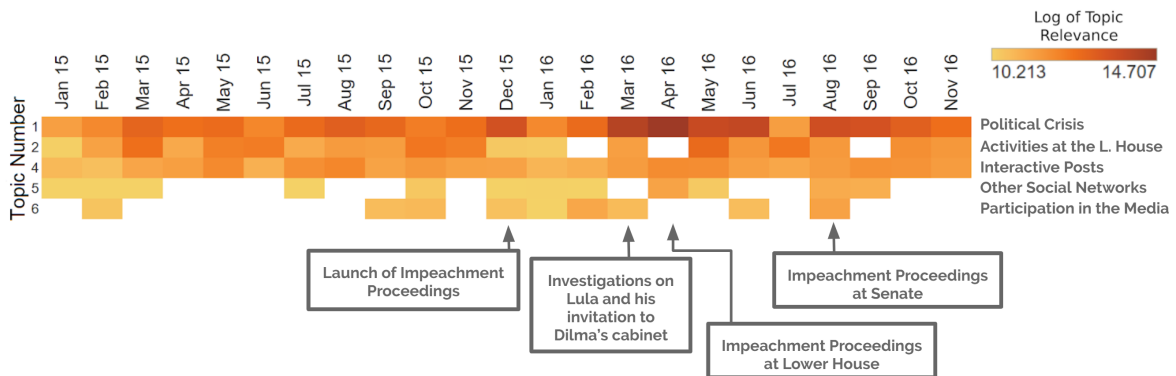


Figure 4.6: Relevance of super topics  $T'_p$  over the months.

- Some top words (“congratulations”, “god”, “friends”, “good”) suggest that  $T'_4$  contains interactive tweets, which may be used by the representatives to get in touch with their friends and/or public.
- Words such as “facebook”, “photos” and “posted” show that  $T'_5$  may cover tweets regarding the participation of the politicians on other social networks.
- $T'_6$  may contain tweets about the participation of politicians on the news media,

Table 4.1: Words describing the super-topics for politicians.

Topic Name	Final Topic	Top 10 Words
Political Crisis	$T'_1$	dilma ( <i>Dilma Rousseff</i> ), brazil, impeachment, government, against, lula ( <i>ex-president Luiz Inacio Lula da Silva</i> ), coup, cunha (Lower House former president, <i>Eduardo Cunha</i> ), temer (vice-president <i>Michel Temer</i> ), today
Activities at the Lower House	$T'_2$	house, committee, representative, today, bill, representatives, meeting, audience, law, minister
Interactive Posts	$T'_4$	day, good, today, congratulations, god, friends, years, week, life, brazil
Other Social Networks	$T'_5$	facebook, posted, new, photo, photos, today, album, mayor, representative, visit
Participation in the Media	$T'_6$	house, today, representative, talk, show, now, TV, day, federal, live

due to the presence of words such as “TV”, “talk” and “show”.

The heat map in Figure 4.6 shows that  $T'_1$  (political crisis) and  $T'_4$  (interactive posts) were discussed all over the period covered by the dataset.  $T'_1$  was more intensely explored than the other topics, especially during December of 2015 and the periods that go from March until June of 2016 and from August to September of 2016. It is also important to notice that tweets about  $T'_2$  (activities at the Lower House) are posted almost over all the time interval covered by the dataset, except for three months (February, April and September of 2016). For the latter case, we can notice that April and September of 2016, when topic  $T'_2$  was not discussed, coincide with the months in which there was the voting of the impeachment at the Lower House (April of 2016) and at the Senate, with the permanent removal of Dilma from the office (last day of August of 2016, start of September), as explained in Section 1.1. At the same time, one can observe that the topic  $T'_1$  (political crisis) was more deeply discussed during these months, as we can see by the more intense colors in the heatmap. This may be the reason why the politicians did not post about their activities, giving more emphasis to the discussion of issues about the political crisis: either they were participating in the voting sessions of the impeachment proceedings (April of 2016) or there was the final voting of the impeachment in the Senate (last day of August of 2016).

Since  $T'_1$  was more deeply discussed, we also calculated the word relevance over time for this topic, in order to better understand specific aspects of the posts. Figure 4.7

shows the 7 most relevant words of this topic and how their relevance changed over the months<sup>2</sup>. Words are ordered by their relevance: the closer to the top of y axis, the more relevant is that word to the topic as a whole.

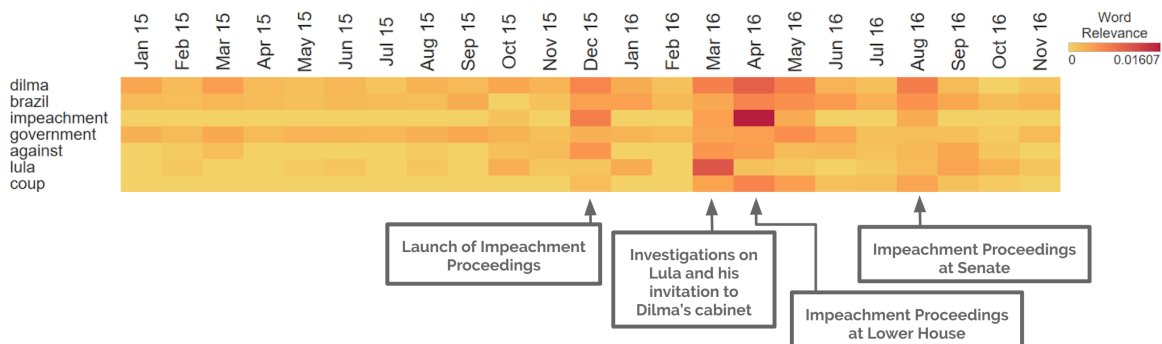


Figure 4.7: Word Relevance over the months - Super Topic 1 (Political Crisis).

As can be observed on the heat map, “dilma” is the most popular word of  $T'_1$  considering all the studied period, which suggests that the ex-president Dilma Rousseff was one of the key points of discussion on the topic concerning the political crisis. Words “Brazil” and “government” are also frequently used over time, possibly indicating discussions about decisions from the government and problems in the country. Despite their high relevance for the topic as a whole (intense colors in the heat map), “impeachment”, “against”, “lula” (*ex-president Lula*) and “coup” record higher values on specific months: December of 2015, March of 2016 and April of 2016. The latter words are probably more relevant in these months due to some remarkable facts, as described on Section 1.1: the launch of impeachment proceedings on December 2015; investigations on Lula and his invitation to join Dilma’s cabinet as the Chief of Staff on March of 2016; and the voting on the impeachment of Dilma at the House of Representatives on April of 2016.

## 4.2 Calculating Polarization from Roll-Call Voting Data

In order to measure polarization among Brazilian politicians, we explore the data from their voting sessions at the Lower House (Section 4.2.1). In our methodology, we proposed a metric to quantify polarity of each representative along the months and we

<sup>2</sup>The complete interactive heat map containing all words in the aggregated topic can be seen on [http://homepages.dcc.ufmg.br/%7Eroberta.coeli/mestrado/WordsRelevanceOverTime\\_SuperTopic1.html](http://homepages.dcc.ufmg.br/%7Eroberta.coeli/mestrado/WordsRelevanceOverTime_SuperTopic1.html)



used the calculated polarity values to compute the polarization index for each month (Section 4.2.2). Our results showed that the polarization for the politicians was higher in 2016 as compared to the previous year. We also observed that the variations to the polarization were mostly related to changes in the distance between gravity centers, i.e., the level of divergence between the opinions of the opposite groups. The experimental results are presented in Section 4.2.3.

### 4.2.1 Data

In this study, we collected roll-call voting data from Brazil’s House of Representatives<sup>3</sup>. This dataset contains all proposed bills, parties orientation and votes from each representative of the House over the years of 2015 and 2016.

We only took into account voting events where Congress leaders of the *Workers Party* (PT) and the *Brazilian Social Democracy Party* (PSDB) gave divergent orientations to their representatives. Since impeached president Dilma Rousseff is a member of PT, Workers Party is a key point in our study. As the major opposition party, PSDB is also relevant to understand how other members of the House changed their support to the ideas of each of these opposite parties along time. The rationale behind this methodology is that, if we include voting events about common-interest subjects, it would not be possible to see the existing ideological differences between the main parties and how the whole set of politicians behave around these contrasting ideas. That is the reason why we included the restriction to only consider the voting sessions that had divergent orientations from the leaders of the main opposite parties.

The final dataset includes 225 voting events that took place between March 10 2015 and December 14 2016, that represents 63.2% of the voting events in which the leaders of the the opposite parties gave a positive or negative orientation. The dataset only comprises the votes of members that have voted “yes”, “no” or “abstention” and have participated of sessions in at least 80% of the studied period (i.e., 16 months), which includes 471 representatives (about 91.8% of the total members of the House).

### 4.2.2 Methods

One of the goals of this work is to measure polarization of Brazilian politicians and investigate its changes over time, for what we propose the following methodology. First,

---

<sup>3</sup>Data is made available by Brazil’s House of Representatives through web services: <http://www.camara.leg.br/SitCamaraWS/Proposicoes.aspx/ListarProposicoesVotadasEmPlenario?ano={0}&tipo=> (proposed bills) and <http://www.camara.leg.br/SitCamaraWS/Proposicoes.aspx/ObterProposicaoPorID?IdProp={0}> (votes).

since we intended to see the temporal evolution of polarization, it was necessary to split our dataset of voting events into time slices. Again, we adopted the month as the time unit. Figure 4.8 shows an overview of the steps of our methodology. Following our approach, we measure polarity  $p_v$  for each representative included in each of the time-sliced sets (**step 1**) and, after that, we computed the Probability Density Function (PDF) for these polarity values of each month (**step 2**). Polarity values and their derived PDFs are used to finally calculate polarization index  $\mu$  across the period (**step 3**).

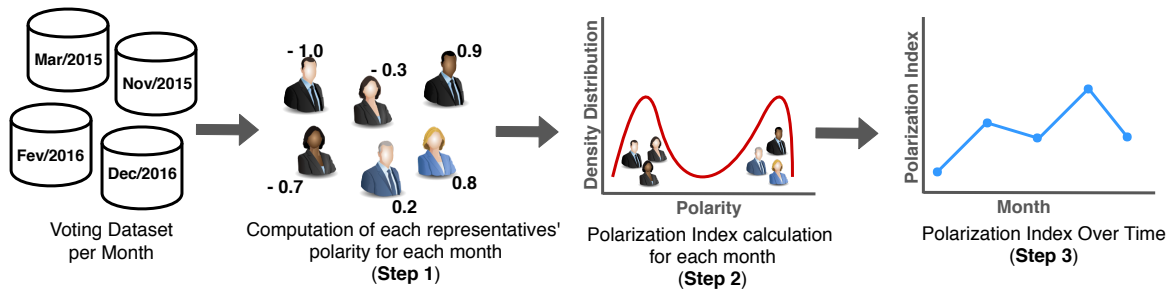


Figure 4.8: Overview of the steps to calculate politicians polarization.

As explained in Chapter 3, we define polarity of an individual as her position around a subject, which she expresses by either agreeing (“yes”), disagreeing (“no”) or being neutral towards it. In this study, polarity is defined as the extension to which a politician agrees with the target party orientation. To compute its value, we took into consideration the votes of each representative in the following situations: i) votes that agreed with the target party position,  $v_{\text{pro}}$ ; ii) votes that were contrary to the target party orientation,  $v_{\text{anti}}$ ; iii) abstention votes,  $v_{\text{abstention}}$ . Workers Party (PT) is taken as the target group due to the fact that it was the ruling party during the period covered by the dataset, to which belongs the impeached president Dilma Rousseff.

First, we compute the total of votes of an individual  $v_T$  as the sum of her votes pro-PT, anti-PT and abstentions (Equation 4.6). We then compute the proportion  $r_i$  of her votes by dividing the number of votes  $v_i$  of each position  $i$  (pro or anti) by the total number of votes of the individual  $v_T$  (Equation 4.7). Having calculated these values for each representative, polarity  $p_v$  can be computed as the difference between the proportion of votes that are similar to the target party orientation ( $r_{\text{pro}}$ ) and the proportion of votes that diverges from it ( $r_{\text{anti}}$ ), as is shown in Equation 4.8.

$$v_T = v_{\text{pro}} + v_{\text{anti}} + v_{\text{abstention}} \quad (4.6)$$

$$r_i = \frac{v_i}{v_T} \quad (4.7)$$

$$p_v = r_{\text{pro}} - r_{\text{anti}} \quad (4.8)$$

Polarity  $p_v$  lies in the range  $[-1, 1]$  and it represents the level of inclination of an individual towards the target party positioning. Here it is important to point out that we decided to reverse the signal of the obtained polarity values  $p_v$ , i.e., the closer the polarity measure is to  $-1.0$ , the more similar are the opinions of the representative to PT; on the other hand, the closer to  $+1.0$ , the more divergent are his opinions to the target party. We chose to associate a negative polarity as the reference of PT positioning – a measure of  $-1.0$  indicates that a representative completely agrees to the party’s opinion – solely due to its left-wing orientation, allowing us to generate more intuitive visualizations.

### 4.2.3 Experimental Results

We began by characterizing the variation of polarities of the Brazilian parties through time. To do this, we calculated the average polarity for each party from the individual polarities of its politicians. Figure 4.9 illustrates the average polarity for all the parties in four months of the dataset: April of 2015, August of 2015, April of 2016 and August of 2016.

As the Figure 4.9 shows, the parties are more evenly distributed over the range of polarities in the months of April and August of 2015. However, in April of 2016, we observe that the average polarities of the parties concentrate in the extremities ( $-1.0$  and  $1.0$ ). For that month, left-wing parties – PT, PSOL, PC do B and REDE – get closer to each other, which means that, on average, members of these parties took decisions that were more similar to PT’s orientation. At the same time, some central and right-wing parties also approximate to PSDB, recording an average polarity close to  $1.0$ . The scenario is similar in August of 2016, when the opposed parties are also concentrated in the extremities, although the central and right-wing parties are more dispersed than on April of 2016.

After the characterization, the individual polarity values were used to compute probability density functions (PDF) for each time slice. Hence, PDFs translate the distribution of opinions from the whole set of representatives across the months, as shown in Figure 4.10. It is important to point out that some months are not shown either because there were no voting session for that month (House recess) or due to the

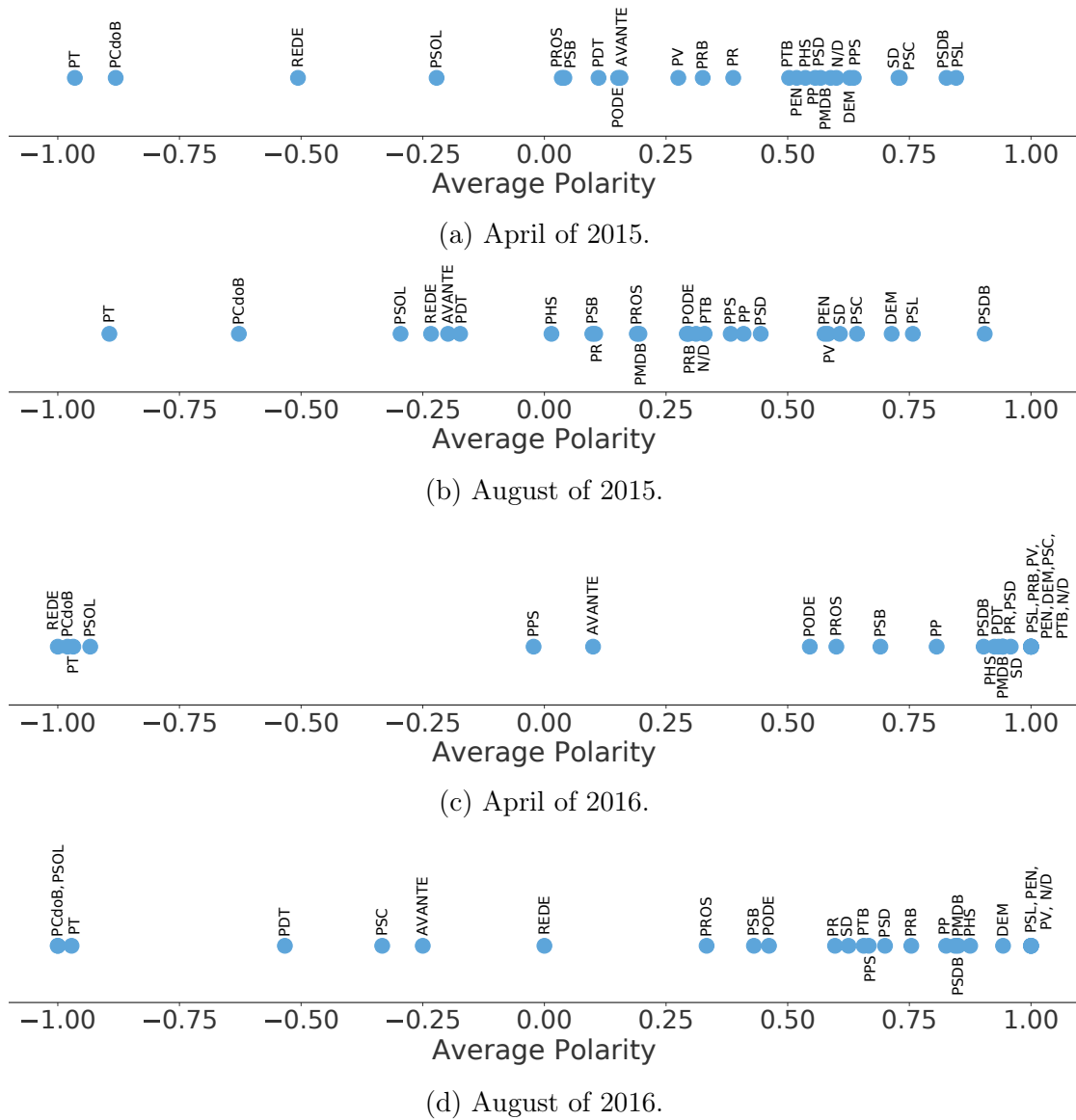


Figure 4.9: Average polarity of the Brazilian parties over time.

restrictions of our study, which only takes into account voting events where PT and PSDB disagree.

Based on these calculated probability density functions, we can observe that:

- Polarity values are more evenly distributed until November of 2015, when representatives seem to start concentrating in opposite groups. In December of the same year, density has considerably increased on the left side, which is almost 3 times greater than the right one. In other words, there was greater support from representatives to the negative polarity (pro-PT).

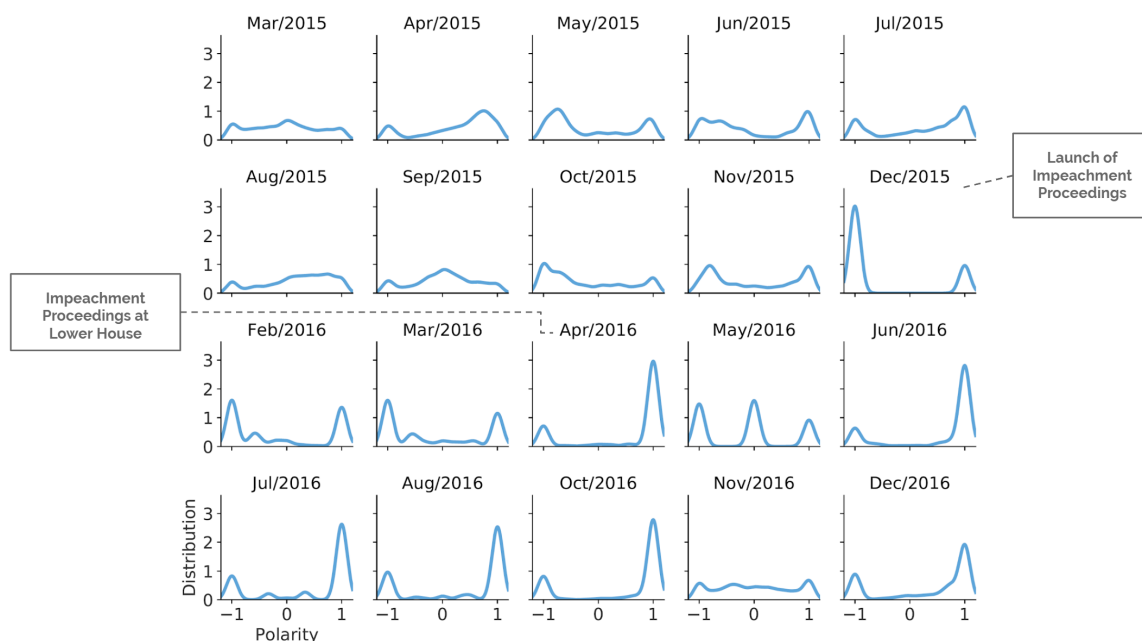


Figure 4.10: Probability density functions for the polarity values of politicians per month.

- In February and March of 2016, the two largest groups that have about the same density are situated in opposite polarities.
- From April of 2016 onwards, we notice that the right-wing group (anti-PT) attracts more representatives, showing a greater density than the pro-PT side. The only exceptions to this situation are May and November of 2016. In May, apart from the opposite groups, there is also a well delimited central group. It possibly indicates that, in this month, some representatives agreed with PT in about half of the voting sessions and were contrary to the party in the other half, which resulted in a polarity value close to 0. In November, polarities are distributed over the whole interval  $[-1, 1]$ , which suggests that most of the representatives did not specifically express support to anti or pro-PT ideas during this month.

In short, one can observe groups of divergent polarities over almost all the studied period, but their differences in the number of representatives are more noticeable after November of 2015. Comparing the whole period, we realize that polarities are more evenly distributed during 2015, whereas there are more noticeable groups of opposite polarities during 2016. We can observe that this change in the behavior of the politicians coincides with the launch of the impeachment proceedings in the beginning of December of 2015, as described in Section 1.1. In 2016, after the proceedings started, the density distributions show two well-defined groups for most of the months, which

seem to indicate that the launch of the impeachment was an important issue that divided the representatives. In April of 2016, we can notice that the right-wing group (anti-PT) is more dense, which also happens in the following months. This situation can be explained by the voting of the impeachment at the House of Representatives in April of 2016, when 367 out of 513 representatives voted to remove Dilma Rousseff from the office, i.e., most of the politicians had a different opinion from her party (PT) position. The similar anti-PT behavior in the subsequent months could be related not only to the impeachment proceedings, but also to the weakening of the Workers Party per se in the view of the corruption scandals and the protests from the population, which also may have affected the political coalitions among the politicians.

As explained in Section 4.2.2, the probability density functions were used to calculate polarization index  $\mu$  for each month. We compare  $\mu$  over the months to understand how polarization changed over the period. Difference in population sizes  $\Delta A$  and distance between gravity centers  $d$  are also compared so as to find aspects related to changes in polarization. Figure 4.11 shows changes in each of these variables over time.

Observing the values of polarization  $\mu$  and its related variables, we can make the following observations:

- Before November of 2015, there were no polarization index values higher than 0.5, whereas  $\mu$  was close to 0.6 for most of the months of 2016.
- In December of 2015, the polarization index reached the highest value for the whole studied period. In the same month, the distance between gravity centers also reached its highest value, which suggests that the increase in polarization is more related to variations in this factor at this point.
- For the year of 2016, the lowest values of polarization index were found to occur in May and November. In both months, there was a decrease in the distance between gravity centers, i.e., the average polarity of each group moved away from each negative or positive end.
- Despite the fact that there are variations to  $\Delta A$  over the entire period, one can notice that the shape of the polarization index resembles the fluctuations to the distance  $d$ . Only on April of 2015, when  $\Delta A$  reaches its peak, we can realize that an increase in this variable causes a major decrease in the polarization index. However, for the most part of the period, variations to  $\Delta A$  seem to have a minor or little influence over polarization when compared to the distance between gravity centers.

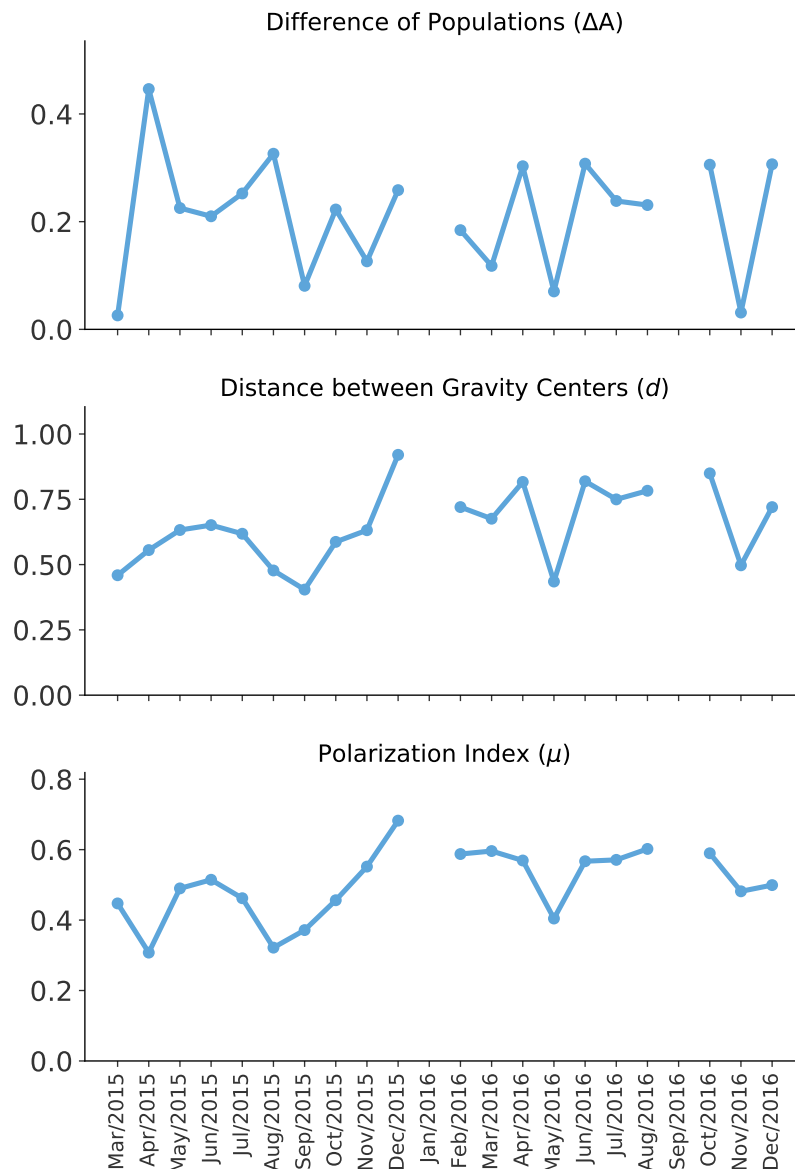


Figure 4.11: Time evolution of polarization index (c), and its related variables: difference in population sizes (a) and distance between gravity centers (b). Some months are not shown either because there was no voting session for that month (Congress recess) or due to the restrictions of our study, which only takes into account voting events that PT and PSDB disagree.

In summary, the polarization index recorded higher values in 2016 as compared to the previous year. In addition, major changes in its value were mostly related to variations in the distance between gravity centers. This situation indicates that, in our case study, polarization among Brazilian representatives was more affected by the average polarity of each group rather than the volume of politicians inside them.

As previously mentioned for the density distributions, the polarization for the

politicians increased after December of 2015, which coincides with the launch of the impeachment proceedings in the start of that month. The polarization values were higher in 2016 than in the previous year, which corroborates our previous observations about the facts around this phenomenon: the representatives became more polarized after the launch of the impeachment proceedings in the end of 2015 and this scenario persisted for the most part of the year of 2016.

In addition, the fact that the polarization among the politicians was more influenced by the distance between the gravity centers indicate that the impeachment and its related events may have resulted in changes in political alliances among the representatives. To put it differently, politicians that had more ideas in common with PT – in the left part of the political spectrum – may have gotten even closer to the party after the events, making similar decisions at the Lower House. On the other hand, politicians with a different ideology – most of which were in the right part of the political spectrum – may have adopted a more opposed attitude, making an opposite decision from PT at the House.

### 4.3 Correlations

Our previous experiments analyzed Twitter messages and roll-call votes from representatives so as to study their behavior in virtual and real-world data. To study associations between these scenarios, we investigated correlations between the polarization index among politicians in voting events and the frequency of their topics in social media, aiming to understand if their actions in the House of Representatives affect what they say on Twitter, and vice versa. To that end, we compared the polarization index  $\mu$  in Congress to the percentage of Twitter topics 1 (Political Crisis) and 2 (Activities at the Lower House) along the months. These topics were chosen because they are the most relevant ones during the covered period, as seen in Section 4.1.

For the purpose of comparing these two variables, Figures 4.12a and 4.12b show the time evolution of the polarization index and the percentage of topics 1 and 2, respectively. Here the percentage of a topic is computed as the number of tweets related to that topic divided by the total of tweets in the month. Besides this qualitative analysis through data visualization, we calculated the Spearman correlation coefficient  $\rho$  for both cases, which results are summarized in Table 4.2.

By observing the correlation results, we cannot reject the hypothesis that there is no association between frequency of tweets from topic 1 and polarization index (p-value = 0.34). By comparison, the corresponding visualization (Figure 4.12a) shows



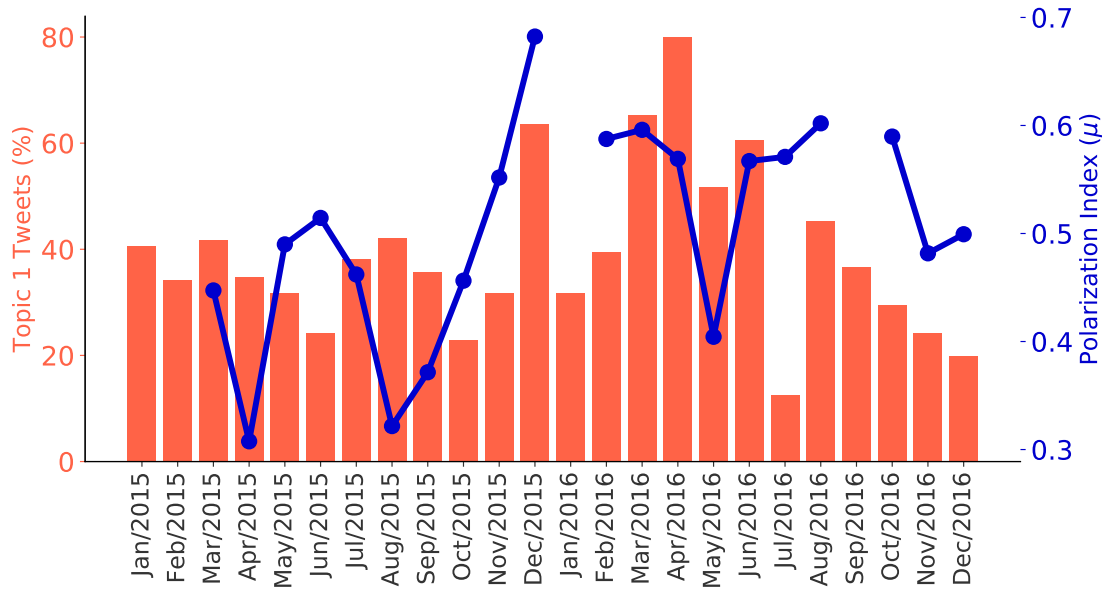
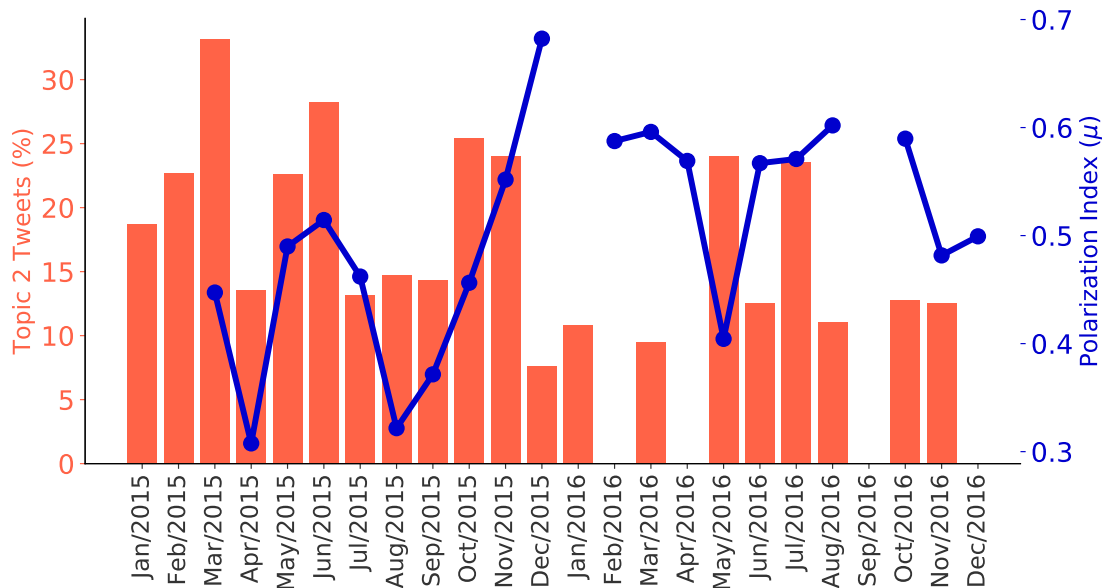
(a)  $T_1'$  (Political Crisis).(b)  $T_2'$  (Activities at the Lower House).

Figure 4.12: Time Evolution of Politicians Polarization Index and Percentage of Tweets of a Topic.

that, between October of 2015 and June of 2016, the polarization index and the topic percentage have a similar growing behavior, even though these variables do not appear to be related in the other months. On the other hand, the correlation between topic 2 and polarization index recorded a small p-value ( $p\text{-value} = 0.03$ ), which suggests that

Topic	Topic Name	Spearman Coefficient ( $\rho$ )	P-Value
$T'_1$	Political Crisis	0.22	0.34
$T'_2$	Activities at the Lower House	-0.53	0.03

Table 4.2: Spearman Correlation Coefficient between Politicians' Percentage of Topics and Polarization Index.

we can reject the null hypothesis. It has a negative coefficient ( $\rho = -0.53$ ), indicating that polarization tends to decrease when topic 2 is more discussed on Twitter. The comparative visualization (Figure 4.12b) shows that there seems to exist a negative correlation between these variables from October of 2015 to June of 2016. However, one can realize that there is a positive relationship between them from March of 2015 to August of 2015, i.e., the polarization index increases as the percentage of tweets from topic 2 grows.

To sum up, it is possible that there is an association between the polarization of the representatives in the Congress and what they say on social media. Nonetheless, we were not able to investigate how they affect each other precisely through the time using the data in our study, because we had only a few observations to compare, not to mention that our voting data does not cover all the months, such as January of 2016 and September of 2016.

# Chapter 5

## Polarization of People

The political crisis in Brazil, combined with recession and the revelation of corruption scandals, such as “Lava Jato” (*Car Wash*), brought thousands of people to the streets to demand the resignation of Dilma. The situation also mobilized protests from supporters of the president to express their disagreement with her removal from office. Not only limited to the streets, anti and pro-Dilma groups also spread their ideas through social networks, which amplified the conflict and gave place to heated political discussions and arguments between the users.

The opening of impeachment proceedings in December of 2015 increased the number of real and virtual manifestations, showing that Brazilians might be becoming more polarized over politics than ever. In this context, Section 5.1 presents our methods and experimental results of evaluating polarization from the Brazilian people on social media. In addition, Section 5.2 brings a comparison of the results from the study of polarization among politicians and the polarization among the Brazilian general public.

### 5.1 Calculating Polarization in Social Media

In order to measure polarization among the general public concerning Dilma’s impeachment process, we collected Twitter messages by using keywords that were related to the most discussed subjects in the Brazilian political scenario (Section 5.1.1). Our methods involved measuring the polarity of users by building a retweet network and, from these values, calculating the polarization index of this set of individuals (Section 5.1.2). Our results showed that the general public recorded high polarization values during the entire studied period. Even though the polarization had small changes over the months, these fluctuations were more related to variations in the difference in the size of the opposite groups, i.e., the polarization among people was mostly affected by the number

of individuals which joined each of the opposite groups. The experimental results are described in Section 5.1.3.

It is important to point out that our methodology for this study of the people polarization is different from the analysis of the politicians in some aspects. First, as the data collection for the general public was based on keywords, performing a topic analysis for this dataset did not seem to bring new information, since the results would be biased to the topics discussed by the politicians. Another point is that, since there are a large number of users in the people dataset, we were able to build a retweet network to calculate their polarization in social media. However, as the number of politicians is smaller, their retweet network comprised only a small percentage of individuals, which is the reason why we only evaluate the polarization of the politicians in their voting dataset.

### 5.1.1 Data

The Brazilian general public tweets were collected through the public Twitter Stream API using the 33 keywords showed in Table 5.1 in the period that goes from March of 2016 to December of 2016. This set of terms was chosen because it includes the main subjects discussed in a preliminary analysis of the politicians topics and were related or had contributed to the Brazilian political instability. For this reason, it comprises names of politicians (`dilma`, `temer`, `lula` and `cunha`, for example), corruption scandals (`lava jato`), companies that were cited for corruption (e.g., `odebrecht`, `petrobras`, `andrade gutierrez`), among other people and institutions that were involved in that turbulent political scenario. The dataset includes approximately 3.3 million users and about 80.4 million tweets that were posted between March 09 2016 and December 27 2016.

Table 5.1: Keywords used to collect tweets for the general public.

<p>lava jato, dilma, impeachment, temer, cunha, odebrecht, moro, lula, petralha, coxinha, renan calheiros, golpe, bolsonaro, globo, camargo correa, andrade gutierrez, empreiteiras, petrobras, delcidio, queiroz galvão, engevix, mendes junior, youssef, collor, romero juca, aecio, anastasia, pizzolatti, paulo roberto costa, renan, delação, delator, policia federal</p>
---

As a preprocessing step, all tweets were lower-cased and stop words were eliminated. In addition, social bots were removed by finding users who posted a large number of messages. We identified that users with more than 700K tweets had aspects of spammers and, thus, these users were removed from the dataset. It is also important

to point out that, despite the fact that our Twitter dataset is in Portuguese, our results were translated to English for the sake of understanding.

### 5.1.2 Methods

The evaluation of polarization among the general public works in a similar fashion of the politicians study, having two basic parts: i) the measurement of individuals polarity, and ii) the computation of polarization index using the previously calculated polarities. For the first part, the measurement of the individuals polarity, we follow the steps shown in Figure 5.1. In this study, we split the dataset into time slices, for which we also consider the month as the time unit (**step 1**). Next, we label a sample of users which posted a list of hashtags that clearly show an opinion (**step 2**). Having a set of labeled users, we build a retweet network (**step 3**), which contains only edges which endpoint is an unlabeled user. Based on the connections between unlabeled and labeled users, we calculate the polarity of each unlabeled user (**step 3**). Having these values, the second part of the methodology involves calculating the probability density functions from the polarity measures. At the end, the final polarization measure is derived from the PDFs.

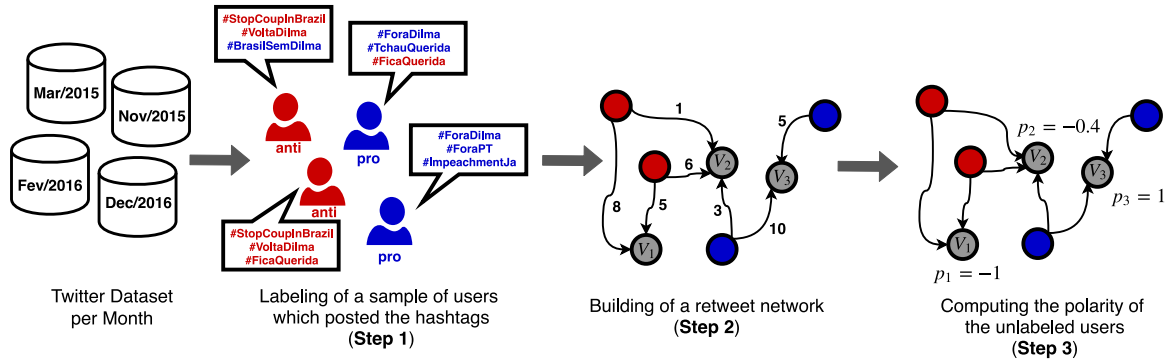


Figure 5.1: Overview of the steps to calculate the polarities of the general public.

As a starting point of our approach, we measure the polarity of each user by finding groups of users with opposite opinions in every time-sliced dataset. Bear in mind that our initial concept of polarity is here instantiated as an user position regarding the impeachment of Dilma Rouseff, i.e., polarity quantifies how much an user supports the resignation of the president.

With the purpose of investigating the different points of view, we began by finding hashtags that clearly show an opinion about the impeachment. To achieve this, the dataset was first characterized so as to find the most popular hashtags related to the

studied subject. Once listed, these hashtags were separated into two groups of opposite opinions using background knowledge from the event, as shown on Table 5.2.

Table 5.2: Hashtags related to the impeachment event in the general public dataset.

<b>Pro-Impeachment Hashtags</b>
#ForaDilma ( <i>Dilma Out</i> ),
#ImpeachmentJa ( <i>Impeachment Now</i> ),
#TchauQuerida ( <i>Goodbye, Dear</i> ),
#DilmaMentirosa ( <i>Dilma is a Liar</i> ),
#ForaPT ( <i>PT Out</i> ),
#BrasilReprovaDilma ( <i>Brazil disapproves Dilma</i> ),
#BrasilSemDilma ( <i>Brazil without Dilma</i> ),
#SenadoVoteSim ( <i>Senate, vote "yes"</i> )
<b>Anti-Impeachment Hashtags</b>
#StopCoupInBrazil ( <i>Stop Coup in Brazil</i> ),
#SOSCoupInBrazil ( <i>SOS, Coup in Brazil</i> ),
#RespeiteAsUrnas ( <i>Respect the Votes</i> ),
#OcupaTudoContraOGolpe ( <i>Occupy Everything Against the Coup</i> ),
#VoltaDilma ( <i>Dilma, Come Back</i> ),
#FicaQuerida ( <i>Stay, Dear</i> ),
#DilmaEInocente ( <i>Dilma is Innocent</i> ),
#SenadoVoteNao ( <i>Senate, vote "no"</i> )

Given a set of users  $U$  in the dataset, we first select from  $U$  the users who posted the listed hashtags. For each of these users, we assign pro/con labels based on the value  $l_u$ , calculated by Equation 5.1, where  $|M_x|$  is the number of messages  $M$  posted by user  $u$  containing hashtags that represent position  $x$ . If  $l_u > 0$ , user  $u$  is labeled as being pro-impeachment, and if  $l_u < 0$ , user  $u$  holds an anti-impeachment position. If  $l_u = 0$ , we are not able to infer the position of the user and she is not labeled.

$$l_u = |M_{\text{pro}}| - |M_{\text{con}}| \quad (5.1)$$

At the end of this process,  $U$  is divided into two subsets:  $U_{\text{labeled}}$  and  $U_{\text{unlabeled}}$ , where the opinion of labeled users is already known. These labeled users  $U_{\text{labeled}}$  were sampled in two equally sized groups of each position. Each group contains 39,940 users, which gives a total of 79,880 users in the sample (about 2.3% of the individuals of the dataset).

In order to calculate polarity for the non-labeled users  $U_{\text{unlabeled}}$ , our method builds a retweet network, for each month slice, by connecting users that retweeted each other from the whole dataset. This network is represented by a weighted directed graph  $G_R = \{V_R, E_R\}$ , where vertices  $V_R$  represent users and edges  $E_R$  connect users  $u_i$  and  $u_j$  ( $u_i \rightarrow u_j$ ) if  $u_j$  retweets a post from  $u_i$ . Besides, edges are weighted by the total

number of retweets. In this network, we only take into account messages from users in  $U_{\text{labeled}}$  that were retweeted by users in  $U_{\text{unlabeled}}$ , i.e., the graph  $G_R$  only contains edges  $U_{\text{labeled}} \rightarrow U_{\text{unlabeled}}$ . All other edges are ignored and all disconnected vertices are also removed. At the end, the final graphs for each time slice included 674,318 non-labeled users, which comprises 22% of the number of users of the whole dataset.

In  $G_R$ , given that an unlabeled user  $u$  in  $U_{\text{unlabeled}}$  retweeted to a set of  $n$  labeled users, this user  $u$  is the endpoint of  $n$  edges  $E_u = \{e_1, e_2, \dots, e_n\}$  in the graph. Each edge  $e_i$  pointing to user  $u$  has a weight  $w_i$ , which represents the number of retweets that user  $u$  made in posts by the connected labeled user. In this way, we can compute the total of retweeted messages  $M_u$  for unlabeled user  $u$  as proposed in Equation 5.2.

$$M_u = \sum_{e_i \in E_u} w_i \quad (5.2)$$

The edges  $E_u$  of user  $u$  can be divided into two groups: a group of edges connecting user  $u$  to labeled users that hold a pro-impeachment position ( $E_u^+$ ); and another group that connects  $u$  to anti-impeachment labeled users ( $E_u^-$ ). Thus, we can determine the total of retweets  $M_u^x$  that user  $u$  made on posts by labeled users of each position  $x$  – pro-impeachment (+) or anti-impeachment (–) – by using Equation 5.3. By using these values, we can also compute the proportion of retweets  $r_u^x$  from unlabeled user  $u$  on labeled users of position  $x$  (Equation 5.4).

$$M_u^x = \sum_{e_i \in E_u^x} w_i \quad (5.3)$$

$$r_u^x = \frac{M_u^x}{M_u} \quad (5.4)$$

Having these values, the polarity  $p_u$  of user  $u$  in  $U_{\text{unlabeled}}$  is calculated according to Equation 5.5. It is defined as the difference in the percentual of retweets from  $u$  on posts of labeled users of each position. Hence, polarity lies in the range  $[-1, 1]$  and represents the level of inclination of an user towards a certain opinion. The closer to 1 it is, the more the user is inclined to a pro-impeachment view. The closer to -1, the more is she inclined to hold an anti-impeachment position.

$$p_u = r_u^+ - r_u^- \quad (5.5)$$

Probability density functions (PDF) were calculated from these polarity measures for each month slice, in order to understand how users are distributed over the different polarities. The derived PDFs were finally used to calculate polarization index  $\mu$ , as

seen on Chapter 3.

### 5.1.3 Experimental Results

As explained in Section 5.1.2, probability density functions were calculated from individual polarity values for each month. They show how users from the dataset are distributed over the polarities, revealing if the individuals are more or less concentrated into groups of some position. Figure 5.2 shows PDFs for every month in the dataset. Observing these distributions, one can notice that most of the individuals are clearly concentrated on divergent groups, with just a few number of users having polarity values close to 0. The left group (anti-impeachment position) has a greater density of individuals for most of the months, except for March of 2016, when the opposite groups seem to comprise about the same number of users.

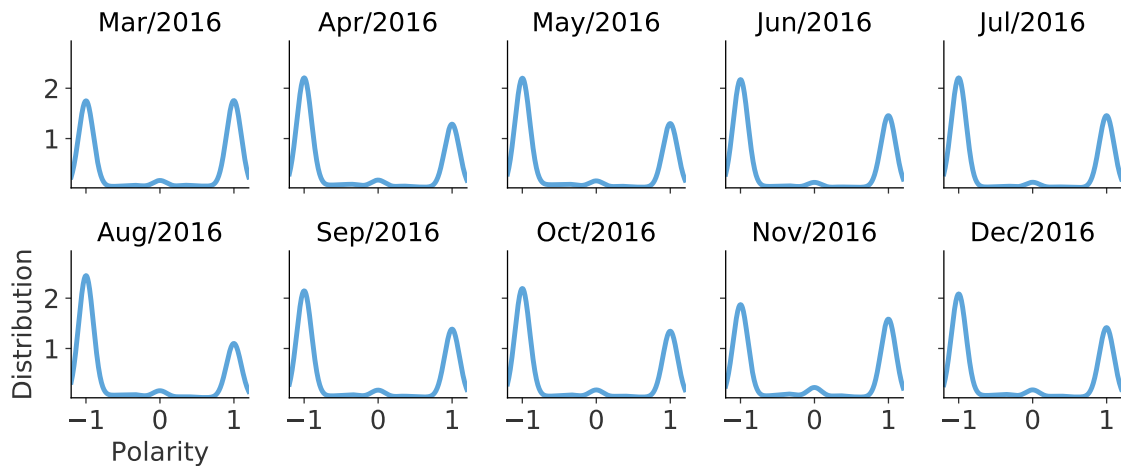


Figure 5.2: Probability density functions for polarity values of Brazilian general public per month.

Once probability density functions were determined, we computed the polarization index  $\mu$  of the general public for each month. Figure 5.3 shows the time evolution of polarization index  $\mu$  and its related variables: difference of populations  $\Delta A$  and distance between gravity centers  $d$ . In the light of these results, we can make the following observations:

- The polarization index has its peak in March ( $\mu = 0.79$ ), when  $\Delta A$  is close to 0. It indicates that groups of opposite opinions had about the same density of users by that time, so that polarity value is largely determined by the distance between the gravity centers.



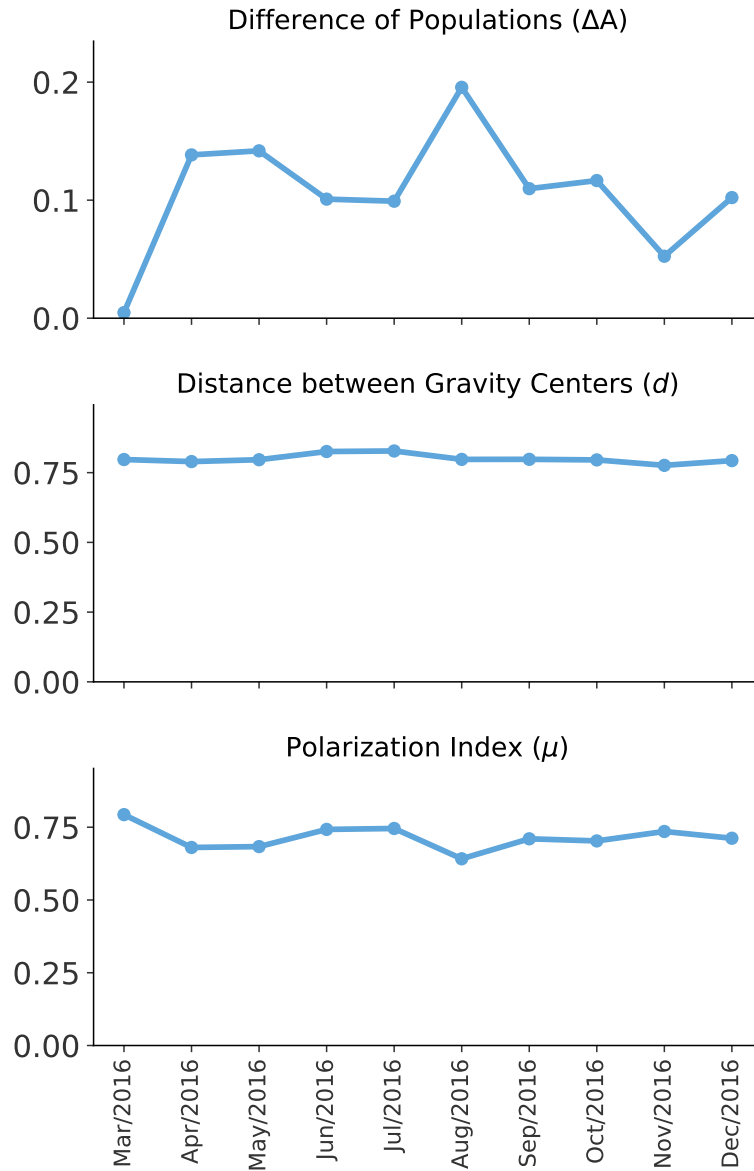


Figure 5.3: Time evolution of polarization index  $\mu$  and its related variables for general public analysis.

- In August, the polarization index recorded its minimum value ( $\mu = 0.64$ ). Since there was no noticeable change in the difference of gravity centers for the adjacent months, we can assume that the decrease in polarization value is mostly related to an increase in  $\Delta A$ , which means that one of the groups became larger than the other. As corroborated by Figure 5.2, the left group has a larger density than the right one in August, and the difference in their sizes appears to be the largest for the whole period.

In short, our results show that the general public recorded high values of polarization over the entire studied period. Since the distance between gravity centers ( $d$ ) remains almost uniform over time, polarization index was mostly affected by fluctuations in the difference in the size of populations ( $\Delta A$ ). In other words, the average position of opposite groups (gravity centers) had little or no variation along the period of study, so that temporal changes in polarization were, for the most part, associated to changes in the number of individuals that are part of each group.

Although the polarization among the general public had small changes along the months, its high values may reflect the tensions among the Brazilian population during the entire studied period, which was characterized by a number of pro and anti-government protests which took place before and after the impeachment proceedings, as seen in Section 1.1. Besides, we can also notice that the highest polarization value was observed in March of 2016, the same month in which the largest anti-government protest in the history of country took place. By the same month, there were also demonstrations of support for the president Dilma, which demands were opposite from the anti-government group. Hence, the highest value in March of 2016 shows that the polarization the population in the social media may be a reflection of the real protests and ideological conflicts in the country.

## 5.2 Comparing People and Politicians Polarization

Recall that one of our main goals is contrasting virtual and real-world data so as to investigate possible associations between them, as it was done by calculating correlations between social media and voting data of Brazilian politicians (Section 4.3). In this study, we investigate these associations by comparing the polarization of politicians – evaluated on their real-world voting data – and the polarization of Brazilian people on social media. To achieve this, we pointed out similarities and differences between the polarization index  $\mu$  and its variables across the months for both cases. Here we intend to understand if the behavior of the general public is related to the actions of the representatives in the Lower House, and vice versa.

We study changes in these variables through a qualitative analysis, since there are only a few number of observations to be compared, which are not enough for a correlation study. Figure 5.4 shows how polarization index  $\mu$  and its related variables – difference between populations sizes  $\Delta A$  and distance between gravity centers  $d$  – changed over time for both politicians and people studies. In order to support our understanding of the visualizations, we also present a statistical summary of each of

these variables in Table 5.3, where we determine the arithmetic mean (Mean), standard deviation (SD), relative standard deviation (RSD), minimum (Min) and maximum (Max) values. It is important to highlight that, although the representatives' dataset covers a broader period range, we calculated the statistics and made our following considerations based on the results from March to December of 2016, since the people dataset is limited to that period.

As can be observed, people recorded high polarization index values for the entire period, having a mean of 0.71; whereas politicians recorded a mean value of 0.54 for the same time interval. The maximum value of polarization index for the politicians ( $\mu = 0.60$ ) is smaller than its minimum value for people ( $\mu = 0.64$ ), which indicates that Brazilian representatives were less polarized than the general public in the studied period. The polarization index of people seems to remain almost uniform across the period, having a small relative standard deviation (RSD = 5.6%). By comparison, politicians polarization index has a larger variation (RSD = 13.0%), especially due to a major decrease in its values in May.

With regard to people results, the small fluctuations in polarization index values seem to be mostly related to the difference between the sizes of their positive and negative populations ( $\Delta A$ ), whereas the distance  $d$  does not have large variations over the months (RSD = 2.5%). For instance, when people record its minimum polarization in August of 2016,  $\Delta A$  reaches its maximum value, while there are no noticeable variations to the distance  $d$  by the same time interval.

The polarization results for the politicians, on the other hand, exhibit more fluctuations over the period as compared to the people. Although there are variations to both  $\Delta A$  and  $d$ , major changes to polarization index were mostly related to fluctuations to the distance  $d$  (see Section 4.2).

Based on these considerations, we observe that the polarization of people is mainly influenced by the difference between the size of populations  $\Delta A$ , whereas changes to polarization of politicians is mostly affected by the distance between the gravity centers  $d$ . It indicates that, for the general public, the polarization process is more related to the number of individuals that join or leave each of the populations. In contrast, the polarization process for Brazilian representatives is more impacted by the level of divergence of the opinions, that is, how different the central opinions of each opposite populations are.

As explained earlier, there were only a few number of observations, in terms of time periods, to be compared in this study due to the restrictions of our datasets. For this reason, there was not enough data to compute correlations. However, our qualitative study was able to point out the differences between the polarization variables

	$\Delta A$		$d$		$\mu$	
	People	Politicians	People	Politicians	People	Politicians
<b>Mean</b>	0.11	0.21	0.80	0.71	0.71	0.54
<b>SD</b>	0.05	0.11	0.02	0.15	0.04	0.07
<b>RSD</b>	45.4%	52.4%	2.5%	21.1%	5.6%	13.0%
<b>Min</b>	0.01	0.03	0.78	0.44	0.64	0.40
<b>Max</b>	0.20	0.31	0.83	0.85	0.79	0.60

Table 5.3: Statistical summary for the difference between populations sizes ( $\Delta A$ ), distance between gravity centers ( $d$ ) and polarity index ( $\mu$ ) in people and politicians studies.

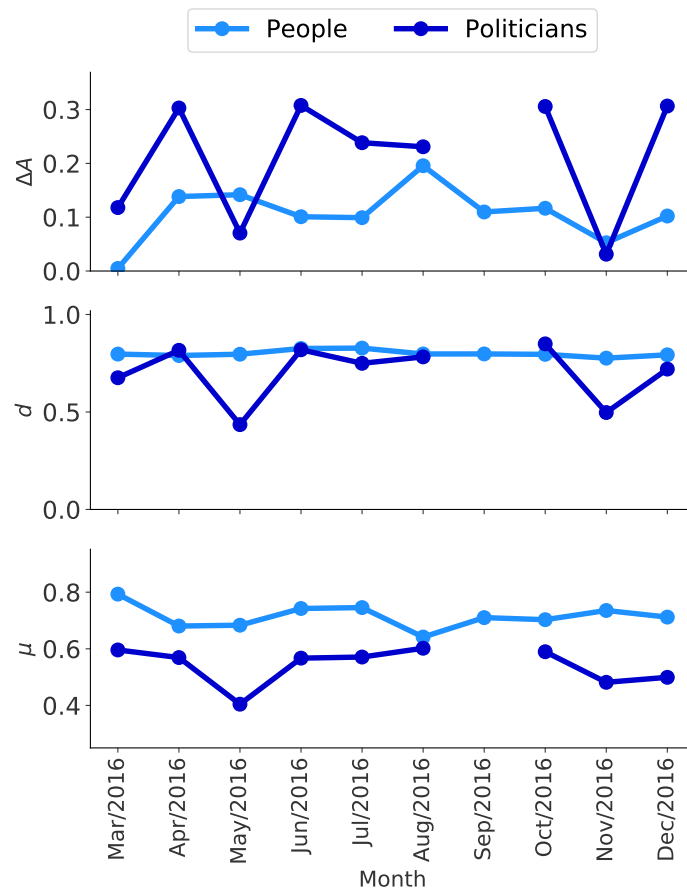


Table 5.4: Comparison between polarization index  $\mu$ , difference between populations sizes  $\Delta A$  and distance between gravity centers  $d$  for people and politicians.

for each case, also revealing that the polarization process was more intense among the general public rather than among Brazilian politicians.

## Chapter 6

# Conclusions and Future Work

The revelation of corruption scandals, such as the *Car Wash* operation, and the impeachment proceedings of Dilma Rousseff in 2016 revealed an intense and profound political and social crisis in Brazil. Tensions emerged not only among Brazilian political parties, but also among common citizens, which took the streets either to demand the impeachment of the president or to demonstrate support for her. Not only limited to the streets, these conflicts were also present in social media platforms, which set the stage for heated discussions among politicians and the general public.

Considering this context, our work aimed at developing computational methods to study political polarization, analyzing online and offline data regarding the impeachment proceedings of 2016 in Brazil. Our study comprised an analysis of the polarization phenomenon among the Brazilian politicians and the general public in the real-world and social media datasets, for which we proposed methods that are able to meet our research objectives and deal with limitations of our data.

With the purpose of understanding the behavior of the politicians in social media, we presented a method to investigate the temporal topic evolution of the Brazilian representatives on Twitter, analyzing their main discussed subjects over time. Our approach divides the studied period into monthly slices and compares the topic models from the different time intervals by building a topic similarity graph. To deal with the sparsity of the tweets, we use the Biterm Topic Model (BTM), an state of art algorithm to find topic models in short texts. We also explored the BTM resulting probabilities to calculate the relevance of the topics and their words over time. As the general public Twitter data collection was based on keywords, we did not use our temporal topic evolution method for the people dataset, since the results of performing such analysis would be biased towards the topics discussed by the politicians.

For the investigation of the temporal topic evolution for the politicians, our ex-

perimental results showed that topics regarding the political crisis and the activities at the Lower House were the most relevant topics discussed by the representatives during the entire period. We also observed that the political crisis topic was more intensely explored than the others, especially after December of 2015, which coincides with the launch of the impeachment proceedings in the Congress. An analysis of the word relevance for the latter topic also revealed that words with a higher relevance value in a specific month were related to events which happened in that month in the impeachment timeline: whereas “dilma”, “Brazil” and “government” were frequently used over time, words such as “impeachment”, “against”, “lula” and “coup” recorded higher relevance for the months of December of 2015, March of 2016 and April of 2016.

In order to quantify polarization for the sets of individuals, we used the polarization metric proposed by [Morales et al. \[2015\]](#), which takes into account the size of the opposite groups and the distance between their central opinions. Since this metric is based on the probability density distribution of the polarities, we needed to measure the opinions of each individual in both the people and the politicians dataset. For the general public, we labeled a sample of users from Twitter and used a retweet network to find the opinions of the unlabeled users. For the politicians, we calculated the polarity of a representative according to how similar was his voting behavior to the Workers Party (PT) orientation. Since the retweet network for the representatives is small and has sparse connections, we solely quantified their polarities and the overall polarization using their voting dataset. Similarly, we do not have the real-world data for the general public, which is the reason why we solely evaluate their polarity and the overall polarization using the social media data.

With regard to the polarization among politicians, our results showed that the representatives recorded higher polarization values in 2016 as compared to the previous year, and these values increased after December of 2015. According to the impeachment timeline, this situation coincides with the launch of the impeachment proceedings in that month, which possibly indicates that the representatives became more polarized after the impeachment started. We also observed that the fluctuations to their polarization values were mostly related to changes in the distance between the gravity centers, i.e., the level of divergence between the central opinions of the opposite groups. This result may also indicate that politicians in the left part of the political spectrum may have gotten even closer to PT after the events, whereas politicians with a different ideology from PT – most of which were in the right part of the political spectrum – may have adopted a more opposed attitude in face of the impeachment.

Concerning the general public, our results showed that people recorded high polarization values over the entire period, having small changes in its value along the

months. This result may reflect the real ideological conflicts among the Brazilian population, meaning that they may have occurred during the whole studied period. In March of 2016, when the polarization reached its highest value in our results, the impeachment timeline points to the occurrence of the largest anti-government protest in the history of Brazil, as well as for demonstrations of support for the president Dilma in the country. These conflicting protests indicate that our polarization results may be a reflection of the ideological differences in the real world.

Having the polarization index for both the politicians and the general public, we performed two different analysis: first, we analyzed the potential correlations between the politicians votes and what they say in their posts, using the discussed topics versus how they voted in the House. Second, we conducted a qualitative analysis to compare the Brazilian population polarization to the politicians polarization.

Regarding the correlations between the frequency of the most relevant topics and the polarization of the politicians, our quantitative results were inconclusive either due to a high p-value for the political crisis topic, or to a small Spearman correlation value for the topic about the activities of the politicians at the Lower House. Although it is possible that there is an association between the polarization of the representatives and what they discuss in social media, our study was not able to point out the associations between these variables precisely.

Since there were only a few observations to the comparison between the politicians polarization and people polarization, we conducted a qualitative analysis to compare them. Our results showed that the polarization was more intense among the general public ( $\bar{\mu} = 0.71$ ) rather than the representatives ( $\bar{\mu} = 0.54$ ). Besides, we also observed that, although the polarization values presented small variations for the general public, these changes were mostly related to the difference between the size of the populations. On the other hand, the polarization among the politicians seems to be more influenced by the distance between the gravity centers. These findings show that, whereas the polarization of people is more affected by the number of individuals that are concentrated in each of the opposite groups, the polarization process for the Brazilian representatives is more impacted by the level of divergence of the central opinions of the opposite groups.

## 6.1 Future Works

As a future work, we suggest a more comprehensive study of the polarization among the general public and among the politicians for the previous years. With more data,

we expect to obtain a more complete understanding of the social dynamics behind the polarization process in society, as well as the ability to investigate and quantify the associations between the polarization phenomenon for both the politicians and the population. In order to understand the causes of the polarization process or the changes in these related variables, a causality analysis is also a good direction as a next step in the research.

Some studies can also enrich the results of our work, bringing more details about the aspects of the political polarization in Brazil. Analyzing the gender difference in the proportion of the votes for the impeachment voting event at the Lower House, for example, would help us understand if men and women had a different position regarding the ex-president Dilma Rousseff, the first woman to hold the Brazilian presidency [Fagundes and Mendonça, 2016]. Another interesting study would involve analyzing the social network of the Brazilian politicians so as to understand the profile of their followers, as well as quantifying the degree of influence of these politicians over their opinions.

Furthermore, investigating other aspects of polarization other than politics would help identifying the factors that contribute to the segregation of opinions in the society and how much political polarization is related to these differences. By analyzing other events, it is also possible to better understand the impact of the “echo chambers” in social media. It is especially important in face of recent studies which show that these structures only appear in some situations, whereas the conversation among people of different ideological views continues to happen for other scenarios, suggesting that the users are not completely immersed into filter bubbles [Barberá et al., 2015].

Another important direction for our work is incorporating a quantitative topic coherence method to our temporal topic evolution method. We currently use a qualitative analysis to evaluate the intermediary results of the method and to define the similarity thresholds. Hence, this step would be important to automatize our method, as well as to measure the quality of our results, also allowing us to compare our method to others in the literature.



# Bibliography

- Abramowitz, A. I. and Saunders, K. L. (2008). Is Polarization a Myth? *The Journal of Politics*, 70(2):542--555.
- Adamic, L. A. and Glance, N. (2005). The Political Blogosphere and the 2004 U.S. Election: Divided they Blog. In *Proceedings of the 3rd International Workshop on Link Discovery*, LinkKDD '05, pages 36--43, New York, New York, USA. ACM Press.
- Baldassarri, D. and Gelman, A. (2008). Partisans without Constraint: Political Polarization and Trends in American Public Opinion. *American Journal of Sociology*, 114(2):408--446. ISSN 0002-9602.
- Barberá, P., Jost, J. T., Nagler, J., Tucker, J. A., and Bonneau, R. (2015). Tweeting from left to right: Is online political communication more than an echo chamber? *Psychological science*, 26(10):1531--1542.
- Bedinelli, T. and Martín, M. (2015). Três grupos organizam os atos anti-Dilma, em meio a divergências. Available at: [https://brasil.elpais.com/brasil/2015/03/13/politica/1426285527\\_427203.html](https://brasil.elpais.com/brasil/2015/03/13/politica/1426285527_427203.html). Last access: November 15, 2018.
- Blei, D. M. and Lafferty, J. D. (2006). Dynamic Topic Models. In *Proceedings of the 23rd International Conference on Machine Learning*, ICML '06, pages 113--120, New York, NY, USA. ACM.
- Blei, D. M., Ng, A. Y., and Jordan, M. I. (2003). Latent dirichlet allocation. *Journal of Machine Learning Research*, 3(Jan):993--1022.
- Bramson, A., Grim, P., Singer, D. J., Fisher, S., Berger, W., Sack, G., and Flocken, C. (2016). Disambiguation of social polarization concepts and measures. *The Journal of Mathematical Sociology*, 40(2):80--111.
- Chen, X., Vorvoreanu, M., and Madhavan, K. P. (2014). Mining Social Media Data for Understanding Students' Learning Experiences. *IEEE Transactions on Learning Technologies*, 7(3):246--259. ISSN 1939-1382.

- Cheng, X., Yan, X., Lan, Y., and Guo, J. (2014). BTM: Topic modeling over short texts. *IEEE Transactions on Knowledge and Data Engineering*, 26(12):2928–2941. ISSN 10414347.
- Cioffi-Revilla, C. (2013). *Introduction to Computational Social Science: Principles and Applications*. Springer Science & Business Media.
- Conover, M., Ratkiewicz, J., Francisco, M. R., Gonçalves, B., Menczer, F., and Flammini, A. (2011). Political polarization on twitter. In *Proceedings of the Fifth International AAAI on Weblogs and Social Media*, volume 133 of *ICWSM '11*, pages 89–96, Barcelona, Spain. Association for the Advancement of Artificial Intelligence.
- Cota, W., Ferreira, S. C., Pastor-Satorras, R., and Starnini, M. (2019). Quantifying echo chamber effects in information spreading over political communication networks. arXiv:1901.03688.
- Das, A., Gollapudi, S., and Munagala, K. (2014). Modeling opinion dynamics in social networks. In *Proceedings of the 7th ACM International Conference on Web Search and Data Mining, WSDM '14*, pages 403–412, New York, NY, USA. ACM.
- Davidov, D., Tsur, O., and Rappoport, A. (2010). Enhanced sentiment learning using twitter hashtags and smileys. In *Proceedings of the 23rd International Conference on Computational Linguistics: Posters, COLING '10*, pages 241–249, Stroudsburg, PA, USA. Association for Computational Linguistics.
- De Choudhury, M., Kiciman, E., Dredze, M., Coppersmith, G., and Kumar, M. (2016). Discovering Shifts to Suicidal Ideation from Mental Health Content in Social Media. In *Proceedings of the 2016 CHI Conference on Human Factors in Computing Systems, CHI '16*, pages 2098–2110, San Jose, CA, USA.
- Diao, Q., Jiang, J., Zhu, F., and Lim, E.-P. (2012). Finding Bursty Topics from Microblogs. In *Proceedings of the 50th Annual Meeting of the Association for Computational Linguistics*, pages 536–544. Association for Computational Linguistics.
- DiMaggio, P., Evans, J., and Bryson, B. (1996). Have american’s social attitudes become more polarized? *American journal of Sociology*, 102(3):690–755.
- Dixit, A. K. and Weibull, J. W. (2007). Political Polarization. *Proceedings of the National Academy of Sciences of the United States of America*, 104(2):7351–7356. ISSN 0027-8424.

- Druckman, J. N., Peterson, E., and Slothuus, R. (2013). How Elite Partisan Polarization Affects Public Opinion Formation. *American Political Science Review*, 107(1):57--79. ISSN 0003-0554.
- Ebrahimi, J., Dou, D., and Lowd, D. (2016a). A joint sentiment-target-stance model for stance classification in tweets. In *Proceedings of COLING 2016, the 26th International Conference on Computational Linguistics: Technical Papers*, pages 2656--2665, Osaka, Japan. The COLING 2016 Organizing Committee.
- Ebrahimi, J., Dou, D., and Lowd, D. (2016b). Weakly supervised tweet stance classification by relational bootstrapping. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 1012--1017.
- Esteban, J.-M. and Ray, D. (1994). On the measurement of polarization. *Econometrica: Journal of the Econometric Society*, pages 819--851.
- Fagundez, I. and Mendonça, R. (2016). Ter 'presidenta' fez diferença para as mulheres? Available at: <https://www.bbc.com/portuguese/brasil-36384962>. Last access: May 01, 2019.
- Farrell, H. (2012). The Consequences of the Internet for Politics. *Annual Review of Political Science*, 15:35--52.
- Farrell, H. and Drezner, D. W. (2008). The power and politics of blogs. *Public Choice*, 134(1-2):15--30.
- Fiorina, M. P. and Abrams, S. J. (2008). Political Polarization in the American Public. *Annual Review of Political Science*, 11(1):563--588. ISSN 1094-2939.
- Garimella, K., Morales, G. D. F., Gionis, A., and Mathioudakis, M. (2018). Quantifying controversy on social media. *ACM Transactions on Social Computing*, 1(1):1--27. ISSN 2469-7818.
- Garimella, K., Weber, I., and De Choudhury, M. (2016). Quote RTs on Twitter: Usage of the New Feature for Political Discourse. In *Proceedings of the 8th ACM Conference on Web Science, WebSci '16*, pages 200--204, New York, New York, USA. ACM, ACM Press.
- Garrett, R. K. (2009). Echo chambers online?: Politically motivated selective exposure among Internet news users. *Journal of Computer-Mediated Communication*, 14(2):265--285.

- Gerbaudo, P. (2018). *Tweets and the streets: Social media and contemporary activism*. Pluto Press.
- Guerra, P. H. C., Jr, W. M., Cardie, C., and Kleinberg, R. (2013). A Measure of Polarization on Social Media Networks Based on Community Boundaries. In *Proceedings of the Seventh International AAAI Conference on Weblogs and Social Media, ICWSM '13*, pages 1--10. Association for the Advancement of Artificial Intelligence.
- IBOPE (2016). Internet e Política: Ativismo nas Redes Sociais. Available at: <http://www.ibopeinteligencia.com/noticias-e-pesquisas/metade-dos-eleitores-brasileiros-receberam-informacoes-sobre-politica-pelo-facebook-twitter-ou-whatsapp>. Last access: November 15, 2018.
- Ileri, I. and Karagoz, P. (2016). Detecting user emotions in twitter through collective classification. In *Proceedings of the 8th International Joint Conference on Knowledge Discovery, Knowledge Engineering and Knowledge Management*, volume 1 of *KDIR '16*, pages 205--212. Science and Technology Publications.
- Jiang, W. and Wu, J. (2017). Active opinion-formation in online social networks. In *IEEE Conference on Computer Communications*, pages 1--9. IEEE.
- Kiritchenko, S., Zhu, X., and Mohammad, S. M. (2014). Sentiment analysis of short informal texts. *Journal of Artificial Intelligence Research*, 50:723--762.
- Kouloumpis, E., Wilson, T., and Moore, J. (2011). Twitter sentiment analysis: The good the bad and the omg! In *Proceedings of the Fifth International AAAI Conference on Weblogs and Social Media, ICWSM '11*, pages 538--541. AAAI.
- Lazer, D., Pentland, A. S., Adamic, L., Aral, S., Barabasi, A. L., Brewer, D., Christakis, N., Contractor, N., Fowler, J., Gutmann, M., et al. (2009). Computational social science. *Science*, 323(5915):721--723.
- Li, J., Li, X., and Zhu, B. (2016). User opinion classification in social media: A global consistency maximization approach. *Information & Management*, 53(8):987 -- 996. ISSN 0378-7206.
- Lietz, H., Wagner, C., Bleier, A., and Strohmaier, M. (2014). When politicians talk: Assessing online conversational practices of political parties on twitter. In *Proceedings of the Eighth International AAAI Conference on Weblogs and Social Media, ICWSM '14*, pages 285--294. AAAI.

- Livne, A., Simmons, M., Adar, E., and Adamic, L. (2011). The party is over here: Structure and content in the 2010 election. In *Proceedings of the Fifth International AAI Conference on Weblogs and Social Media, ICWSM '11*, page SI, Barcelona, Spain. AAI.
- Mei, Q. and Zhai, C. (2005). Discovering evolutionary theme patterns from text: an exploration of temporal text mining. In *Proceedings of the eleventh ACM SIGKDD international conference on Knowledge discovery in data mining*, pages 198--207. ACM.
- Mohammad, S., Kiritchenko, S., Sobhani, P., Zhu, X., and Cherry, C. (2016). Semeval-2016 task 6: Detecting stance in tweets. In *Proceedings of the 10th International Workshop on Semantic Evaluation (SemEval-2016)*, pages 31--41.
- Mohammad, S. M. (2012). #emotional tweets. In *Proceedings of the First Joint Conference on Lexical and Computational Semantics, SemEval '12*, pages 246--255, Stroudsburg, PA, USA. Association for Computational Linguistics.
- Mohammad, S. M., Sobhani, P., and Kiritchenko, S. (2017). Stance and sentiment in tweets. *ACM Transactions Internet Technology*, 17(3):1--23. ISSN 1533-5399.
- Morales, A. J., Borondo, J., Losada, J. C., and Benito, R. M. (2015). Measuring political polarization: Twitter shows the two sides of Venezuela. *Chaos: An Interdisciplinary Journal of Nonlinear Science*, 25(3):033114. ISSN 1089-7682.
- Morinaga, S. and Yamanishi, K. (2004). Tracking dynamics of topic trends using a finite mixture model. In *Proceedings of the tenth ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 811--816. ACM.
- Newman, M. E. J. (2006). Modularity and community structure in networks. *Proceedings of the National Academy of Sciences*, 103(23):8577--8582. ISSN 0027-8424.
- O'Connor, B., Balasubramanyan, R., Routledge, B. R., Smith, N. A., et al. (2010). From tweets to polls: Linking text sentiment to public opinion time series. In *Proceedings of the Fourth International AAI Conference on Weblogs and Social Media*, volume 4 of *ICWSM '10*, pages 122--129. AAI.
- Pang, B. and Lee, L. (2008). Opinion mining and sentiment analysis. *Foundations and Trends in Information Retrieval*, 2(1-2):1--135.

- Rabelo, J. C. B., Prudêncio, R. B. C., and Barros, F. A. (2012). Using link structure to infer opinions in social networks. In *2012 IEEE International Conference on Systems, Man, and Cybernetics (SMC)*, pages 681–685. ISSN 1062-922X.
- Rajadesingan, A. and Liu, H. (2014). Identifying users with opposing opinions in twitter debates. In *International conference on social computing, behavioral-cultural modeling, and prediction*, pages 153–160. Springer.
- Ravi, K. and Ravi, V. (2015). A survey on opinion mining and sentiment analysis: tasks, approaches and applications. *Knowledge-Based Systems*, 89:14–46.
- Ribeiro, V. A. and Gomes Goveia, F. (2016). A Comissão do Impeachment na Rede: o Histórico das Narrativas Políticas Sobre o Impedimento de Dilma Rousseff no Twitter. In *Anais do XVIII Congresso de Ciências da Comunicação na Região Nordeste*, page S.I. Intercom – Sociedade Brasileira de Estudos Interdisciplinares da Comunicação.
- Ruediger, M. A., Martins, R., da Luz, M., and Grassi, A. (2014). Ação coletiva e polarização na sociedade em rede para uma teoria do conflito no brasil contemporâneo. *Revista Brasileira de Sociologia*, 2(4):205–234.
- Schmitt, J. (2016). How to measure ideological polarization in party systems. In *ECPR Graduate Student Conference*, page SI. University of Tartu.
- Sobhani, P., Mohammad, S., and Kiritchenko, S. (2016). Detecting stance in tweets and analyzing its interaction with sentiment. In *Proceedings of the Fifth Joint Conference on Lexical and Computational Semantics*, pages 159–169.
- Stilo, G. and Velardi, P. (2016). Efficient temporal mining of micro-blog texts and its application to event discovery. *Data Mining and Knowledge Discovery*, 30(2):372–402.
- Sunstein, C. R. (2002). The Law of Group Polarization. *Journal of Political Philosophy*, 10(2):175–195. ISSN 0963-8016.
- Tatagiba, L. (2018). Entre as ruas e as instituições: os protestos e o impeachment de dilma rousseff. *Lusotopie*, 17(1).
- Taulé, M., Martí, M. A., Rangel, F. M., Rosso, P., Bosco, C., Patti, V., et al. (2017). Overview of the task on stance and gender detection in tweets on catalan independence at ibereval 2017. In *2nd Workshop on Evaluation of Human Language*

- Technologies for Iberian Languages, IberEval 2017*, volume 1881, pages 157--177. CEUR-WS.
- Velasco, C., D'Agostino, R., Reis, T., et al. (2016). Da eleição à votação do impeachment. Available at: <http://especiais.g1.globo.com/politica/2016/processo-de-impeachment-de-dilma/da-eleicao-a-votacao-do-impeachment/>. Last access: November 15, 2018.
- Venceslau, P. (2016). Grupos e partidos de oposição se juntam em atos contra dilma. Available at: <https://exame.abril.com.br/brasil/pela-1a-vez-grupos-e-partidos-de-oposicao-se-associam-nos-atos-contradilma/>. Last access: November 15, 2018.
- Wang, X. and McCallum, A. (2006). Topics over time: a non-markov continuous-time model of topical trends. In *Proceedings of the 12th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 424--433. ACM.
- Watts, J. (2014). Dilma rousseff pledges unity after narrow brazil election victory. Available at: <https://www.theguardian.com/world/2014/oct/26/brazil-re-elects-dilma-rousseff-president>. Last access: October 6, 2018.
- Wilson, T., Wiebe, J., and Hoffmann, P. (2005). Recognizing contextual polarity in phrase-level sentiment analysis. In *Proceedings of the conference on human language technology and empirical methods in natural language processing*, pages 347--354. Association for Computational Linguistics.
- Wolfsfeld, G., Segev, E., and Sheaffer, T. (2013). Social Media and the Arab Spring: Politics Comes First. *The International Journal of Press/Politics*, 18(2):115--137.
- Xie, W., Zhu, F., Jiang, J., Lim, E. P., and Wang, K. (2016). TopicSketch: Real-time bursty topic detection from twitter. *IEEE Transactions on Knowledge and Data Engineering*, 28(8):2216--2229. ISSN 10414347.
- Yan, X., Guo, J., Lan, Y., and Cheng, X. (2013). A biterm topic model for short texts. In *Proceedings of the 22nd international conference on World Wide Web, WWW '13*, pages 1445--1456, Rio de Janeiro, Brazil. ACM Press. ISSN 145032035X.
- Yin, H., Cui, B., Lu, H., Huang, Y., and Yao, J. (2013). A unified model for stable and temporal topic detection from social media data. In *2013 IEEE 29th International Conference on Data Engineering (ICDE)*, pages 661--672. IEEE. ISSN 10844627.

Zhang, X., Chen, X., Chen, Y., Wang, S., Li, Z., and Xia, J. (2015). Event detection and popularity prediction in microblogging. *Neurocomputing*, 149:1469--1480.