Quantifying Controversy in Social Media

Kiran Garimella Aalto University Helsinki, Finland kiran.garimella@aalto.fi

Aristides Gionis Aalto University Helsinki, Finland aristides.gionis@aalto.fi

ABSTRACT

Which topics spark the most heated debates in social media? Identifying these topics is a first step towards creating systems which pierce echo chambers. In this paper, we perform a systematic methodological study of controversy detection using social media network structure and content.

Unlike previous work, rather than identifying controversy in a single hand-picked topic and use domain-specific knowledge, we focus on *comparing topics in any domain*. Our approach to quantifying controversy is a graph-based threestage pipeline, which involves (i) building a *conversation* graph about a topic, which represents alignment of opinion among users; (ii) partitioning the conversation graph to identify potential sides of the controversy; and (iii) measuring the amount of controversy from characteristics of the graph.

We perform an extensive comparison of controversy measures, as well as graph building approaches and data sources. We use both controversial and non-controversial topics on Twitter, as well as other external datasets. We find that our new random-walk-based measure outperforms existing ones in capturing the intuitive notion of controversy, and show that content features are vastly less helpful in this task.

1. INTRODUCTION

Given their widespread diffusion, online social media are becoming increasingly important in the study of social phenomena such as peer influence, framing, bias, and controversy. Ultimately, we would like to understand how users perceive the world through the lens of their social media feed. For instance, to offer users the possibility to balance their "news diet" [20, 21] on controversial topics by recommending contrarian content, which supports a view that differs from what they are mostly exposed to [25]. However, before addressing these advanced application scenarios, we first need to focus on the fundamental yet challenging task of distinguishing whether a topic of discussion is controversial, Our work is motivated by interest in observing controversies at societal level, monitoring their evolution, and possibly understanding which issues become controversial and why.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

WSDM 2016, February 22-25, 2016, San Francisco, California, USA

C 2016 Copyright held by the owner/author(s). Publication rights licensed to ACM. ISBN ACM 978-1-4503-3716-8/16/02. DOI: 10.1145/1235

Gianmarco De Francisci Morales Aalto University Helsinki, Finland gdfm@acm.org

Michael Mathioudakis HIIT, Aalto University Helsinki, Finland michael.mathioudakis@hiit.fi

The study of controversy in social media is not new; there are many previous studies aimed at identifying and characterizing controversial issues, mostly around political debates [1, 7, 23, 24] but also for other topics [15]. And while most recent papers have focused on Twitter [7, 15, 23, 24], controversy in other social-media platforms, such as blogs [1] and opinion fora [2], have also been analyzed.

However, most previous papers have severe limitations. First, the majority of previous studies focus on controversy regarding political issues, and in particular, they are centered around long-lasting major events, such as elections [1, 7]. More crucially, most previous works can be characterized as *case studies*, where controversy is identified in a single carefully-curated dataset, collected using ample domain knowledge and auxiliary domain-specific sources (e.g., an extensive list of hashtags regarding a major political event, or a list of left-leaning and right-leaning blogs).

We aim to overcome those limitations. Our goal is to identify controversy regarding topics in any domain (e.g., political, economical, or cultural), and in a general setting, i.e., without prior domain-specific knowledge about the topic in question. In addition, we aim at comparing different topics, in order to find the most controversial ones. These properties allow to deploy a system in-the-wild, and are valuable for building real-world applications.

In order to enable such a versatile framework, we work with topics that are defined in a lightweight and domain-agnostic manner. Specifically, when focusing on Twitter, a topic can be specified as a text query. For example, "#beefban" is a special keyword (a "hashtag") that was employed by Twitter users to signal that their posts referred to a decision of the Indian government, in March 2015, about the consumption of beef meat in India. In this case, the query "#beefban" defines a topic of discussion, and the related activity consists of all posts that contain the query.

We represent a topic of discussion with a *conversation* graph. In such a graph, vertices represent people, and edges represent conversation activity and interactions, such as *posts*, *comments*, *mentions*, or *endorsements*. Our working hypothesis is that it is possible to analyze the conversation graph of a topic to reveal how controversial the topic is. In particular, we expect the conversation graph of a controversial topic to have a *clustered structure*. This hypothesis is based on the fact that a controversial topic entails different sides with opposing points of view, and individuals on the same side tend to endorse and amplify each other's arguments [1, 2, 7]. Our main contribution is to test this hypothesis. We achieve this by studying a large number of candidate features, based on the following *aspects* of activity: (i) structure of endorsements, i.e., who agrees with whom on the topic, (ii) structure of the social network, i.e., who is connected with whom among the participants in the conversation, (iii) content, i.e., the keywords used in the topic, (iv) sentiment, i.e., the tone (positive or negative) used to discuss the topic. Our study shows that, except from content-based features, all the other ones are useful in detecting controversial topics, to different extents. Particularly for Twitter, we find the endorsement features (i.e., retweets) to be the most useful.

The extracted features are then used to compute the *controversy score* of a topic. We offer a systematic definition and provide a thorough evaluation of measures to quantify controversy. We employ a broad range of topics, both controversial and non-controversial ones, on which we evaluate several measures, either defined in this paper or coming from the literature [15, 24]. We find that one of our newly-proposed measure, based on *random walks*, is able to discriminate controversial topics with great accuracy. In addition, it also generalizes well as it agrees with previously-defined measures when tested on datasets from existing work. We also find that the *variance* of the sentiment expressed on a topic is a reliable indication of controversy.

The approach to quantifying controversy presented in this paper can be condensed into a three-stage pipeline: (i) build a *conversation graph* among the users who contribute to a topic, where edges signify that two users are in agreement, (ii) identify the potential sides of the controversy from the graph structure or the textual content, and (iii) quantify the amount of controversy in the graph.

The rest of this paper is organized as follows. Section 2 discusses how this work fills gaps in the existing literature. Subsequently, Section 3 provides a high level description of the pipeline for quantifying controversy for a topic, while Sections 4, 5, and 6 detail each stage. We empirically evaluate the proposed measures of controversy in Section 7. Section 8 extends the evaluation to a few measures that do not fit the pipeline. We conclude in Section 9 with a discussion on possible improvements and directions for future work, as well as lessons learned from carrying out this study.

2. RELATED WORK

Analysis of controversy in online news and social media has attracted considerable attention, and a number of papers have provided very interesting case studies. In one of the first papers, Adamic and Glance [1] study the link patterns and discussion topics of political bloggers, focusing on blog posts on the U.S. presidential election of 2004. They measure the degree of interaction between liberal and conservative blogs, and provide evidence that conservative blogs are linking to each other more frequently and in a denser pattern. These findings are confirmed by the more recent study of Conover et al. [7], who also study controversy in political communication regarding congressional midterm elections. Using data from Twitter, Conover et al. identify a highly segregated partisan structure (present in the retweet graph, but not in the mention graph), with limited connectivity between leftand right-leaning users. In another recent work related to controversy analysis in political discussion, Mejova et al. [23] identify a significant correlation between controversial issues and the use of negative affect and biased language.

The papers mentioned so far study controversy in the political domain, and provide case studies centered around long-lasting major events, such as presidential elections. In this paper, we aim to identify and quantify controversy for any topic discussed in social media, including short-lived and ad-hoc ones. The problem we study has been considered by previous work, but the methods proposed so far are, to a large degree, domain-specific.

The work of Conover et al. [7], discussed above, employs the concept of modularity and graph partitioning in order to verify (but not quantify) controversy structure of graphs extracted from discussion of political issues on Twitter. In a similar setting, Guerra et al. [15] propose an alternative graph-structure measure. Their measure relies on the analysis of the boundary between two (potentially) polarized communities, and performs better than modularity. Differently from these studies, our contribution consists in providing an extensive study of a large number of measures, including the ones proposed earlier, and demonstrating clear improvement over those. We also aim at quantifying controversy in diverse and in-the-wild settings, rather than carefully-curated domain-specific datasets.

In a recent study, Morales et al. [24] quantify polarity via the propagation of opinions of influential users on Twitter. They validate their measure with a case study from Venezuelan politics. Again, our methods are not only more general and domain agnostic, but they provide more intuitive results. In a different approach, Akoglu proposes a polarization metric that uses signed bipartite opinion graphs [2]. The approach differs from ours as it relies on the availability of this particular type of data, which is not as readily available as social-interaction graphs.

Similarly to the papers discussed above, in our work we quantify controversy based on the graph structure of social interactions. In particular, we assume that controversial and polarized topics induce graphs with clustered structure, representing different opinions and points of view. This assumption relies on the concept of "echo chambers," which states that opinions or beliefs stay inside communities created by like-minded people, who reinforce and endorse the opinions of each other. This phenomenon has been quantified in many recent studies [3, 11, 14].

A different direction for quantifying controversy followed by Choi et al. [6] and Mejova et al. [23] relies on text and sentiment analysis. Both studies focus on language found on news articles. In our case, since we are mainly working with Twitter, where text is short and noisy, and since we are aiming at quantifying controversy in a domain-agnostic manner, text analysis has its limitations. Nevertheless, we experiment with incorporating content in our approach.

Finally, our findings on controversy have many potential applications on news-reading and public-debate scenarios. For instance, quantifying controversy can provide a basis for analyzing the "news diet" of readers [20, 21], offering the chance of better information by providing recommendations of contrarian views [25], or trying to deliberate debates [10] and connect people with opposing opinions [9, 13].

3. PIPELINE

Our approach to measuring controversy is based on a systematic way of characterizing social media activity. We employ a pipeline with three stages, namely graph building, graph partitioning, and measuring controversy. The final output of the pipeline is a value that measures how controversial a topic is, with higher values corresponding to higher degree of controversy. We provide a high-level description of each stage here and more details in the sections that follow.

3.1 Building the Graph

The purpose of this stage is to build a *conversation graph* that represents activity related to a single *topic* of discussion. In our pipeline, a topic is operationalized as a *query*, and the social media activity related to the topic consists of those items (e.g., posts) that match the given query. For example, in the context of Twitter, the query might simply consist of a keyword, such as "#ukraine", in which case the related activity consists of all tweets that contain that keyword. Even though we describe textual queries in standard document-retrieval form, in principle queries can take other forms, as long as they are able to induce a graph from the social media activity (e.g., RDF queries, or topic models).

Each item related to a topic is associated with one user who generated it, and we build a graph where each user who contributed to the topic is assigned to one vertex. In this graph, an edge between two vertices represents *endorsment*, *agreement*, or *shared point of view* between the corresponding users. Section 4 details several ways to build such a graph.

3.2 Partitioning the Graph

In the second stage, the resulting conversation graph is fed into a graph partitioning algorithm to extract two partitions (we defer considering multi-sided controversies to a further study). Intuitively, the two partitions correspond to two disjoint sets of users who possibly belong to different sides in the discussion. In other words, the output of this stage answers the following question: "assuming that users are split into two sides according to their point of view on the topic, which are these two sides?" Section 5 describes this stage in further detail. If indeed there are two sides which do not agree with each other –a controversy– then the two partitions should be loosely connected to each other, given the semantic of the edges. This property is captured by a measure computed in the third and final stage of the pipeline.

3.3 Measuring Controversy

The third and last stage takes as input the graph built by the first stage and partitioned by the second stage, and computes the value of a *controversy measure* that characterizes how controversial the topic is. Intuitively, a controversy measure aims to capture how separated the two partitions are. We test several such measures, including ones based on random walks, betweenness centrality, and low-dimensional embeddings. Details are provided in Section 6.

4. GRAPH BUILDING

This section provides details about the different approaches we follow to build graphs from raw data. We use posts on Twitter to create our datasets.¹ Twitter is a natural choice for the problem at hand, as it represents one of the main fora for public debate in online social media, and is often used to report news about current events. Following the procedure described in Section 3.1, we specify a set of queries, and build one graph for each query. The set of topics we choose is balanced between controversial and non-controversial ones, so as to test for both false positives and false negatives.

We use Twitter hashtags as queries. Users employ hashtags to indicate the topic of discussion their posts pertain to. Among the large number of hashtags that appear in the Twitter stream, we consider hashtags that were trending during the period from Feb 27 to Jun 15, 2015. By manual inspection we find that most trending hashtags are not related to controversial discussions.

We first manually pick a set of 10 hashtags that we know represent *controversial* topics of discussion. All topics in this set have been widely covered by mainstream media, and have generated ample discussion both online and offline. Moreover, to have a dataset that is balanced between controversial and non-controversial topics, we sample another set of 10 hashtags that represent *non-controversial* topics of discussion. These hashtags are related mostly to soft news and entertainment, but also to events that, while being impactful and dramatic, did not generate large controversies (e.g., #nepal and #germanwings). In addition to our intuition that these topics are non-controversial, we manually check a sample of tweets and, are unable to identify any controversy.²

For each hashtag, we retrieve all tweets that contain it and are generated during the observation window. We also ensure that the selected hashtags are associated with a large enough volume of activity. Table 1 presents the final set of hashtags, along with their description and the number of related tweets.³ For each hashtag, we build a graph G where we assign a vertex to each user who employs it, and generate edges according to one of the following four approaches.

1. Retweet graph. Typically, retweets are used as endorsements. Users who retweet signal endorsement of the opinion expressed in the original tweet by propagating it further. Retweets are not constrained to occur only between users who are connected in Twitter's social network, but users are allowed to re-post tweets generated by any other user.

We select the edges for graph G based on the retweet activity in the topic: an edge exists between two users u and v if there are at least two ($\tau = 2$) retweets between them that use the hashtag, irrespective of direction. We remark that, in preliminary experimentation with this approach, building the retweet graph with a threshold $\tau = 1$ did not produce reliable results. We presume that a single retweet on a topic is not enough of a signal to infer endorsement. Using $\tau = 2$ retweets as threshold proves to be a good trade-off between high selectivity (which hinders analysis) and noise reduction. The resulting size for each retweet graph is listed in Table 1.

2. Follow graph. In this approach, we build the follow graph induced by a given hashtag. We select the edges for graph G based on the social connections between Twitter users who employ the given hashtag: an edge exists between users u and v if u follows v or vice-versa. We stress that the graph G built with this approach is topic-specific, as the edges in G are constrained to connections between users who discuss the topic that is specified as input to the pipeline.

The rationale for using this graph is based on an assumption of the presence of homophily in the social network, which is a common trait in this setting. To be more precise, we

¹From the full Twitter firehose stream.

 $^{^2{\}rm Code}$ and networks used in this work are available at http://github.com/gvrkiran/controversy-detection.

³We use a hashtag in Russian, #MapIII, which we refer to as #russia_march from here on, for convenience.

Hashtag	# Tweets	Retweet graph		Follow graph		Description and collection period (2015)
		V	E	V	E	
#beefban	84543	1610	1978	799	6026	Government of India bans beef, Mar 2–5
#nemtsov	183477	6546	10172	2156	46529	Death of Boris Nemtsov, Feb 28–Mar 2
#netanyahuspeech	254623	9434	14476	4292	297136	Netanyahu's speech at U.S. Congress, Mar 3–5
#russia_march	118629	2134	2951	1189	16471	Protests after death of Boris Nemtsov ("march"), Mar 1–2
#indiasdaughter	167704	3659	4323	1542	9480	Controversial Indian documentary, Mar 1–5
#baltimoreriots	218157	3902	4505	1441	28291	Riots in Baltimore after police kills a black man, Apr 28–30
#indiana	116379	2467	3143	946	24328	Indiana pizzeria refuses to cater gay wedding, Apr 2–5
#ukraine	287438	5495	9452	3383	84035	Ukraine conflict, Feb 27–Mar 2
#gunsense	318409	7106	11483	1821	103840	Gun violence in U.S., Jun 1–30
#leadersdebate	1139344	25983	44174	9566	344088	Debate during the U.K. national elections, May 3
#sxsw	343652	9304	11003	4558	91356	SXSW conference, Mar 13–22
#1dfamheretostay	501960	15292	26819	3151	20275	Last OneDirection concert, Mar 27–29
#germanwings	907510	29763	39075	2111	7329	Germanwings flight crash, Mar 24–26
#mothersday	1798018	155599	176915	2225	14160	Mother's day, May 8
#nepal	1297995	40579	57544	4242	42833	Nepal earthquake, Apr 26–29
#ultralive	364236	9261	15544	2113	16070	Ultra Music Festival, Mar 18–20
#FF	408326	5401	7646	3899	63672	Follow Friday, Jun 19
#jurassicworld	724782	26407	32515	4395	31802	Jurassic World movie, Jun 12-15
#wcw	156243	10674	11809	3264	23414	Women crush Wednesdays, Jun 17
#nationalkissingday	165172	4638	4816	790	5927	National kissing day, Jun 19

Table 1: Datasets statistics: hashtag, sizes of the follow and retweet graphs, and description of the event. The top group represent controversial topics, while the bottom one represent non-controversial ones.

expect that on a given topic people will agree more often than not with people they follow, and that for a controversial topic of discussion this phenomenon will be reflected in well-separated partitions of the resulting graph. Note that using the entire social graph would not necessarily produce well-separated partitions that correspond to single topics of discussion, as those partitions would be "blurred" by the existence of additional edges that are due to other reasons (e.g., offline social connections).

On the practical side, while the retweet information is readily available in the stream of tweets, the social network of Twitter is not. Collecting the follower graph thus requires an expensive crawling phase. The resulting graph size for each follow graph is listed in Table 1.

3. Content graph. We create the edges of graph G based on whether users post instances of the same content. Specifically, we experiment with the following three variants: create an edge between two vertices if the users (i) use the same hashtag, other than the one that defines the topic, (ii) share a link to the same URL, or (iii) share a link with the same URL domain (e.g., cnn.com is the domain for all pages on the website of CNN).

4. Hybrid content & retweet graph. We create edges for graph G according to a state-of-the-art process that blends content and graph information [26]. Concretely, we associate each user with a vector of frequencies of mentions for different hashtags. Subsequently, we create edges between pairs of users whose corresponding vectors have high cosine similarity, and combine them with edges from the retweet graph, built as described above. For details, we refer the interested reader to the original publication [26].

5. GRAPH PARTITIONING

As previously explained, we use a graph partitioning algorithm to produce two partitions on the conversation graph. To do so, we rely on a state-of-the-art off-the-shelf algorithm, METIS [19]. Figure 1 displays the two partitions returned for some of the topics on their corresponding retweet and follow graphs (Figures 1(a)-(d) and Figures 1(e)-(h), respectively).⁴ The partitions are depicted in blue or red. The graph layout is produced by Gephi's ForceAtlas2 algorithm [16], and is based solely on the structure of the graph, not on the partitioning computed by METIS.

From an initial visual inspection of the partitions identified on retweet and follow graphs, we find that the partitions match well with our intuition of which topics are controversial (the partitions returned by METIS are well separated for controversial topics). To make sure that this initial assessment of the partitions is not an artifact of the visualization algorithm we use, we try other layouts offered by Gephi. In all cases we observe similar patterns. We also manually sample and check tweets from the partitions, to verify the presence of controversy. While this anecdotal evidence is hard to report, indeed the partitions seem to capture the spirit of the controversy.⁵

On the contrary, the partitions identified on content graphs fail to match our intuition. All three variants of the contentbased approach lead to sparse graphs and highly overlapping partitions, even in cases of highly controversial issues. The same pattern applies for the hybrid approach, as shown in Figure 2. We also try a variant of the hybrid graph approach with vectors that represent the frequency of different URL domains mentioned by a user, with no better results. We thus do not consider these approaches to graph building any further in the remainder of this paper.

Finally, we try graph partitioning algorithms of other types. Besides METIS (cut based), we test spectral clustering, label propagation, and affiliation-graph-based models. The difference among these methods is not significant, however from visual inspection METIS generates the cleanest partitions.

⁴Other topics show similar trends, omitted for lack of space. ⁵For instance, of these two tweets for #netanyahuspeech from two users on either side, one is clearly supporting the speech https://t.co/OVeWB4XqIg, while the other highlights



Figure 1: Sample conversation graphs with retweet (top) and follow (bottom) features (visualized using the force-directed layout algorithm in Gephi). The left side is controversial, (a,e) #beefban, (b,f) #russia_march, while the right side is non-controversial, (c,g) #sxsw, (d,h) #germanwings.



Figure 2: Partitions obtained for (a) #beefban, (b) #russia_march by using the hybrid graph building approach. The partitions are more noisy than those in Figures 1(a,b).

6. CONTROVERSY MEASURES

This section describes the controversy measures used in this work. For completeness, we describe both those measures proposed by us (§6.1, 6.2, 6.3) as well as the ones from the literature that we use as baselines (§6.4, 6.5).

6.1 Random walk

This measure uses the notion of random walks on graphs. It is based on the rationale that, in a controversial discussion, there are authoritative users on both sides, as evidenced by a large degree in the graph. The measure captures the intuition of how likely a random user on either side is to be exposed to authoritative content from the opposing side. Let G(V, E) be the graph built by the first stage and its two partitions X and Y, $(X \cup Y = V, X \cap Y = \emptyset)$ identified by the second stage of the pipeline. We first distinguish the *k* highest-degree vertices from each partition. High-degree is a proxy for authoritativeness, as it means that a user has received a large number of endorsements on the specific topic. Subsequently, we select one partition at random (each with probability 0.5) and consider a random walk that starts from a random vertex in that partition. The walk terminates when it visits any high-degree vertex (from either side).

We define the Random Walk Controversy (RWC) measure as follows. "Consider two random walks, one ending in partition X and one ending in partition Y, RWC is the difference of the probabilities of two events: (i) both random walks started from the partition they ended in and (ii) both random walks started in a partition other than the one they ended in." The measure is quantified as

$$RWC = P_{XX}P_{YY} - P_{YX}P_{XY}, \tag{1}$$

where P_{AB} , $A, B \in \{X, Y\}$ is the conditional probability

 $P_{AB} = P(\text{start in partition } A \mid \text{end in partition } B).$

The aforementioned probabilities have the following desirable properties: (i) they are not skewed by the size of each partition, as the random walk starts with equal probability from each partition, and (ii) they are not skewed by the total degree of vertices in each partition, as the probabilities are conditional on ending in either partition (i.e., the fraction of random walks ending in each partition is irrelevant). *RWC* is close to one when the probability of crossing sides is low, and close to zero when the probability of crossing sides is comparable to that of staying on the same side.

the negative reactions https://t.co/v9RdPudrrC.

6.2 Betweenness

Let us consider the set of edges $C \subseteq E$ in the cut defined by the two partitions X, Y. This measure uses the notion of edge betweenness and how the betweenness of the cut differs from that of the other edges. Recall that the betweenness centrality bc(e) of an edge e is defined as

$$bc(e) = \sum_{s \neq t \in V} \frac{\sigma_{s,t}(e)}{\sigma_{s,t}},$$
(2)

where $\sigma_{s,t}$ is the total number of shortest paths between vertices s, t in the graph and $\sigma_{s,t}(e)$ is the number of those shortest paths that include edge e.

The intuition here is that, if the two partitions are wellseparated, then the cut will consist of edges that bridge *structural holes* [5]. In this case, the shortest paths that connect vertices of the two partitions will pass through the edges in the cut, leading to high betweenness values for edges in C. On the other hand, if the two partitions are not well separated, then the cut will consist of *strong ties*. In this case, the paths that connect vertices across the two partitions will pass through one of the many edges in the cut, leading to betweenness values for C similar to the rest of the graph.

Given the distributions of edge betweenness on the cut and the rest of the graph, we compute the KL divergence d_{KL} of the two distributions by using kernel density estimation to compute the PDF and sampling 10 000 points from each of these distributions (with replacement). We define the **Betweenness Centrality Controversy** (*BCC*) measure as

$$BCC = 1 - e^{-d_{KL}},\tag{3}$$

which assumes values close to zero when the divergence is small, and close to one when the divergence is large.

6.3 Embedding

This measure is based on a low-dimensional embedding of graph G produced by Gephi's ForceAtlas2 algorithm [16] (the same algorithm used to produce the plots in Figures 1 and 2). We opt for this algorithm as it produces well-separated plots for controversial topics.

Let us consider the two-dimensional embedding $\phi(v)$ of vertices $v \in V$ produced by ForceAtlas2. Given the partition X, Y produced by the second stage of the pipeline, we calculate the following quantities:

- d_X and d_Y , the average embedded distance among pairs of vertices in the same partition, X and Y respectively,
- d_{XY} , the average embedded distance among pairs of vertices across the two partitions X and Y.

Inpsired by the Davies-Bouldin (DB) index [8], we define the Embedding Controversy measure EC as

$$EC = 1 - \frac{d_X + d_Y}{2d_{XY}}.$$
(4)

EC is close to one for controversial topics, corresponding to better-separated graphs and thus to higher degree of controversy, and close to zero for non-controversial topics.

6.4 Boundary Connectivity

This controversy measure was proposed by Guerra et al. [15], and is based on the notion of boundary and internal vertices. Let $u \in X$ be a vertex in partition X; u belongs to the *boundary* of X iff it is connected to at least one

vertex of the other partition Y, and it is connected to at least one vertex in partition X that is not connected to any vertex of partition Y. Following this definition, let B_X, B_Y be the set of boundary vertices for each partition, and $B = B_X \cup B_Y$ the set of all boundary vertices. By contrast, vertices $I_X = X - B_X$ are said to be the *internal* vertices of partition X (similarly for I_Y). Let $I = I_X \cup I_Y$ be all internal vertices in either partition. The reasoning for this measure is that, if the two partitions represent two sides of a controversy, then boundary vertices will be more strongly connected to internal vertices than to other boundary vertices of either partition. This intuition is captured in the formula

$$GMCK = \frac{1}{|B|} \sum_{u \in B} \frac{d_i(u)}{d_b(u) + d_i(u)} - 0.5$$
(5)

where $d_i(u)$ is the number of edges between vertex u and internal vertices I, while $d_b(u)$ is the number of edges between vertex u and boundary vertices B. Higher values of the measure correspond to higher degrees of controversy.

6.5 Dipole Moment

This controversy measure was presented by Morales et al. [24], and is based on the notion of *dipole moment* that has its origin in physics. Let $R(u) \in [-1, 1]$ be a polarization value assigned to vertex $u \in V$. Intuitively, extreme values of R (close to -1 or 1) correspond to users who belong most clearly to either side of the controversy. To set the values R(u) we follow the process described in the original paper [24]: we set $R = \pm 1$ for the top-5% highest-degree vertices in each partition X and Y, and set the values for the rest of the vertices by label-propagation. Let n^+ and n^- be the number of vertices V with positive and negative polarization values, respectively, and ΔA the absolute difference of their normalized size $\Delta A = \left|\frac{n^+ - n^-}{|V|}\right|$. Moreover, let gc^+ (gc^-) be the average polarization value among vertices n^+ (n^-) and set d as half their absolute difference, $d = \frac{|gc^+ - gc^-|}{2}$. The dipole moment controversy measure is defined as

$$MBLB = (1 - \Delta A)d. \tag{6}$$

The rationale for this measure is that, if the two partitions X and Y are well separated, then label propagation will assign different extreme (±1) *R*-values to the two partitions, leading to higher values of the *MBLB* measure. Note also that larger differences in the size of the two partitions (reflected in the value of ΔA) lead to decreased values for the measure, which takes values between zero and one.

7. EXPERIMENTS

In this section we report the results of the various configurations of the pipeline proposed in this paper. As previously stated, we omit results for the content and hybrid graph building approaches presented in Section 4 as they do not perform well. We instead focus on the retweet and follow graphs, and test all the measures presented in Section 6 on the Twitter topics described in Table 1. In addition, we test all the measures on a set of external datasets used in previous studies [1, 7, 15] to validate the measures against a known ground truth. Finally, we use an evolving dataset from Twitter collected around the death of Venezuelan president Hugo Chavez [24] to show the *evolution* of the controversy measures in response to high-impact events.



Figure 3: Controversy scores on *retweet* graphs of various controversial and non-controversial datasets

To avoid potential overfitting, we use only eight graphs as testbed during the development of the measures, half of them controversial (beefban, nemtsov, netanyahu, russia_march) and half non-controversial (sxsw, germanwings, onedirection, ultralive). This procedure resembles a 40/60% train/test split in traditional machine learning applications.⁶

Twitter hashtags. Figure 3 and Figure 4 report the scores computed by each measure for each of the 20 hashtags, on the retweet and follow graph, respectively. Each figure shows a set of beanplots,⁷ one for each measure. Each beanplot shows the estimated probability density function for a measure computed on the topics, the individual observations are shown as small white lines in a one-dimensional scatter plot, and the median as a longer black line. The beanplot is divided into two groups, one for controversial topics (left/dark) and one for non-controversial ones (right/light). A larger separation of the two distributions indicates that the measure is better at capturing the characteristics of controversial topics. For instance, this separation is fundamental when using the controversy score as a feature in a classification algorithm.

Figures 3 and 4 clearly show that RWC is the best measure on our datasets. BCC and EC show varying degrees of separation and overlap, although EC performs slightly better as the distributions are more concentrated, while BCC has a very wide distribution. The two baselines GMCK and MBLBinstead fail to separate the two groups. Especially on the retweet graph, the two groups are almost indistinguishable.

For all measures the median score of controversial topics is higher than for non-controversial ones. This result suggests that both graph building methods, retweet and follow, are able to capture the difference between controversial and noncontroversial topics. Given the broad range of provenience of the topics covered by the dataset, and their different characteristics, the consistency of results is very encouraging.



⁷A beanplot is an alternative to the boxplot for visual comparison of univariate data among groups.



Figure 4: Controversy scores on *follow* graphs of various controversial and non-controversial datasets.

External datasets. We have shown that our approach works well on a number of datasets extracted in-the-wild from Twitter [12]. However, how well does it generalize to datasets from different domains?

We obtain a comprehensive group of datasets kindly shared by authors of previous works: *Political blogs*, links between blogs discussing politics in the US [1]; *Twitter politics*, Twitter messages pertaining to the 2010 midterm election in US [7]; and the following five graphs used in the study that introduced *GMCK* [15], (a) *Gun control*, retweets about gun control after the shooting at the Sandy Hook school; (b) *Brazil soccer*, retweets about to two popular soccer teams in Brazil; (c) *Karate club*, the well-known social network by Zachary [28]; (d) *Facebook university*, a social graph among students and professors at a Brazilian university; (e) *NYC teams*, retweets about two New York City sports teams.

Table 2 shows a comparison of the controversy measures under study on the aforementioned datasets.⁸ For each dataset we also report whether it was considered controversial in the original paper, which provides a sort of "ground truth" to evaluate the measures against.

All the measures are able to distinguish controversial graphs to some extent, in the sense that they return higher values for the controversial cases. The only exception is Karate club. Both RWC and MBLB report low controversy scores for this graph. It is possible that the graph is too small for such random-walk-based measures to function properly. Conversely, BCC is able to capture the desired behavior, which suggests that shortest-path and random-walk based measures might have a complementary function.

Interestingly, while the Political blogs datasets is often considered a gold standard for polarization and division in online political discussions, all the measures agree that it presents only a moderate level of controversy. On the other hand, the Twitter politics dataset is clearly one of the most

⁸The datasets provided by Guerra et al. [15] are slightly different from the ones used in the original paper because of some irreproducible filtering used by the authors. We use the datasets provided to us verbatim.

Table 2: Results on external datasets. The 'C?' column indicates whether the previous study considered the dataset controversial (ground truth).



Figure 5: Controversy scores on 56 retweet graphs from Morales et al. Day 'D' (indicated by the blue vertical line) indicates the announcement of the death of president Hugo Chavez.

controversial one across all measures. This difference might suggest that the measures are more geared towards capturing the dynamics of controversy as it unfolds on social media, which might differ from more traditional blogs. For instance, one such difference is the *cost* of an endorsement: placing a link on a blog post arguably consumes more mental resources than clicking on the retweet button.

For the Gun control dataset, Guerra et al. needs to manually distinguish three different partitions in the graph: gun rights advocates, gun control supporters, and moderates. Our pipeline is able to find the two communities with opposing views (grouping together gun control supporters and moderates) without any external help. All measures agree with the conclusions drawn in the original paper that this topic is highly controversial.

Evolving controversy. We have shown that our approach also generalizes well to datasets from different domains. But in a real deployment the measures need to be computed continuously, as new data arrives. How well does our method work in such a setting? And how do the controversy measures evolve in response to high-impact events?

To answer these questions, we use a datasets from the study that introduced MBLB [24]. The dataset comprises Twitter messages pertaining to political events in Venezuela around the time of death of Hugo Chavez (Feb-May 2013). The authors built a retweet graph for each of the 56 days (one graph per day) around the day of the death.

Figure 5 shows how the intensity of controversy evolves according to the measures under study (which occurs on day 'D'). The measure proposed in the original paper, *MBLB*, which we use as 'ground truth', shows a clear decrease of controversy on the day of the death, followed by a progressive increase in the controversy of the conversation. The original interpretation states that on the day of the death a large amount of people, also from other countries, retweeted news of the event, creating a single global community that got together at the shock of the news. After the death, the ruling and opposition party entered in a fiery discussion over the next elections, which increased the controversy.

All the measures proposed in this work show the same trend as *MBLB*. Both *RWC* and *EC* follow very closely the original measure (Pearson correlation coefficients r of 0.944 and 0.949, respectively), while *BCC* shows a more jagged behavior in the first half of the plot (r = 0.743), due to the discrete nature of shortest paths. All measures however present a dip on day 'D', an increase in controversy in the second half, and another dip on day 'D+20'. Conversely, *GMCK* reports an almost constant moderate value of controversy during the whole period (r = 0.542), with barely noticeable peaks and dips. We conclude that our measures generalize well also to the case of evolving graphs, and behave as expected in response to high-impact events.

8. CONTENT

In this section we explore alternative approaches to measuring controversy that use only the content of the discussion rather than the structure of user interactions. As such, these methods do not fit in the pipeline described in Section 3. The question we address is "does content help in measuring the controversy of a topic?" In particular, we test two types of features extracted from the content. The first, is a typical IR-inspired bag-of-words representation. The second instead is based on NLP tools for sentiment analysis.

8.1 Content as bag of words

We take in input the raw content of the social media posts, in our case the Tweets containing a specific hashtag. We represent each tweet as a vector in a high-dimensional space composed of the words used in the whole topic, after standard preprocessing used in IR (lowercasing, stopword removal, stemming). Following the lines of our main pipeline, we group these vectors in two clusters by using CLUTO [18] with cosine distance.

The underlying assumption is that the two sides, while sharing the use of the hashtag for the topic, use different vocabularies in reference to the issue at hand. For example, for #beefban a side may be calling for "freedom" while the opposing one for "respect." We use KL divergence as a measure of distance between the vocabularies of the two clusters, and the I2 measure [22] of clustering heterogeneity.

We use an unpaired Wilcoxon rank-sum test at the p = 0.05 significance level, but we are unable to reject the null hypothesis that there is no difference in these measures between the controversial and non-controversial topics. Therefore, there is not enough signal in the content representation to discern between controversial and non-controversial topics with confidence. This result suggests that the bag-of-words representation of content is not a good basis for our task. It also agrees with our earlier attempts to use content to build the graph used in the pipeline (see Section 4) – which suggests that using content for the task of quantifying controversy might not be straightforward.

8.2 Sentiment Analysis

Next, we resort to NLP techniques for sentiment analysis to analyze the content of the discussion. We use SentiStrength [27] trained on tweets to give a sentiment score in [-4, 4] to each tweet for a given topic. In this case we do not try to cluster tweets by their sentiment. Rather, we analyze the difference in distribution of sentiment between controversial and non-controversial topics.

While it is not possible to say that controversial topics are more positive or negative than non-controversial ones (results omitted due to space constraints), we can detect a difference in their variance. Indeed, controversial topics have a higher variance than non-controversial ones, as shown in Figure 6. Controversial ones have a variance of at least 2, while non-controversial ones have a variance of at most 1.5.

In practice, the "tones" with which controversial topics are debated are stronger, and sentiment analysis is able to detect this aspect. While this signal is clear, it is not straightforward to incorporate it into the measures based on graph structure. Moreover, this feature relies on technologies that do not work reliably for languages other than English and hence cannot be applied for topics such as #russia_march.

9. DISCUSSION

The task we tackle in this work is certainly not an easy one, and this study has some limitations, which we discuss in this section. We also report a set of negative results that we produced while coming up with the measures presented. We believe these results will be very useful in steering this research topic towards a fruitful direction.



Figure 6: Sentiment variance controversy score for controversial and non-controversial topics.

9.1 Limitations

Twitter only. We present our findings mostly on datasets coming from Twitter. While this is certainly a limitation, Twitter is one of the main venues for online public discussion, and one of the few for which data is available. Hence, Twitter is a natural choice. In addition, our measures generalize well to datasets from other social media and the Web.

Choice of data. We manually pick the controversial topics in our dataset, which might introduce bias. In our choice we represent a broad set of typical controversial issues coming from religious, societal, racial, and political domains. Unfortunately, ground truths for controversial topics are hard to find, especially for ephemeral issues. However, the topics are unanimously judged controversial by the authors. Moreover, the hashtags represent the intuitive notion of controversy that we strive to capture, so human judgement is an important ingredient we want to use.

Overfitting. While this work presents the largest systematic study on controversy in social media so far, we use only 20 topics for our main experiment. Given the small number of examples, the risk of overfitting our measures to the dataset is real. We reduce this risk by using only 40% of the topics during the development of the measures. Additionally, our measures agree with previous independent results on external datasets, which further decreases the likelihood of overfitting.

Reliance on graph partitioning. Our pipeline relies on a graph partitioning stage, whose quality is fundamental for the proper functioning of the controversy measures. Given that graph partitioning is a hard but well studied problem, we rely on off-the-shelf techniques for this step. A measure that bypasses this step entirely is highly desirable, and we report a few unsuccessful attempts in the next subsection.

Multisided controversies. Not all controversies involve only two sides with opposing views. Some times discussions are multifaceted, or there are three or more competing views on the field. The principles behind our measures neatly generalize to multisided controversies. However, in this case the graph partitioning component needs to automatically find the optimal number of partitions. We defer experimental study of such cases to an extended version of this paper.

9.2 Negative Results

We briefly review a list of methods that failed to produce reliable results and were discarded early in the process of refining our controversy measures.

Mentions graph. Conover et al. [7] rely on the mention graph in Twitter to detect controversies. However, in our dataset the mention graphs are extremely sparse given that we focus on short-lived events. Merging the mentions into the retweet graph does not provide any noticeable improvement. Previous studies have also shown that people retweet similar ideologies but mention across ideologies [4]. We exploit this intuition by using correlation clustering for graph partitioning, with negative edges for mentions. Alas, the results are qualitatively worse than those obtained by METIS.

Cuts. Simple measures such as size of the cut of the partitions do not generalize across different graphs. Conductance (in all its variants) also yields poor results. Prior work identifies controversies by comparing the structure of the graph with randomly permuted ones [7]. Unfortunately, we obtain equally poor results by using the difference in conductance with cuts obtained by METIS and by random partitions.

Community structure. Good community structure in the conversation graph is often understood as a sign that the graph is polarized or controversial. However, this is not always the case. We find that both assortativity and modularity (which have been previously used to identify controversy) do not correlate with the controversy scores, and are not good predictors for how controversial a topic is. The work by Guerra et al. [15] presents clear arguments and examples of why modularity should be avoided.

Partitioning. As already mentioned, bypassing the graph partitioning to compute the measure is desirable. We explore the use of the all pairs expected hitting time computed by using SimRank [17]. We compute the SPID (ratio of variance to mean) of this distribution, however results are mixed.

9.3 Conclusions

In this paper, we performed the first large-scale systematic study for quantifying controversy in social media. We have shown that previously-used measures are not reliable and demonstrated that controversy can be identified both in the retweet and topic-induced follow graph. We have also shown that simple content-based representations do not work in general, while sentiment analysis offers promising results.

Among the measures we studied, the random-walk-based RWC most neatly separates controversial topics from noncontroversial ones. Besides, our measures gracefully generalize to datasets from other domains and previous studies.

This work opens several avenues for future research. First, it is worth exploring alternative approaches and testing additional features, such as, following a generative-model-based approach, or exploiting the temporal evolution of the discussion of a topic. From the application point of view, the controversy score can be used to generate recommendations that foster a healthier "news diet" on social media.

Acknowledgements This work is supported by the European Community's H2020 Program under the scheme "INFRAIA-1-2014-2015: Research Infrastructures," grant agreement #654024 "SoBigData: Social Mining & Big Data Ecosystem." We would like to thank Morales et al. [24] for providing the code and networks for the *MBLB* method.

10. REFERENCES

- L. A. Adamic and N. Glance. The political blogosphere and the 2004 us election: divided they blog. In *LinkKDD*, pages 36–43, 2005.
- [2] L. Akoglu. Quantifying political polarity based on bipartite opinion networks. In *ICWSM*, 2014.
- [3] J. An, D. Quercia, and J. Crowcroft. Partisan sharing: facebook evidence and societal consequences. In COSN, pages 13–24, 2014.

- [4] A. Bessi, G. Caldarelli, M. Del Vicario, A. Scala, and W. Quattrociocchi. Social Determinants of Content Selection in the Age of (Mis)Information. In *Social Informatics*, pages 259–268, 2014.
- [5] R. S. Burt. Structural holes: The social structure of competition. Harvard university press, 2009.
- [6] Y. Choi, Y. Jung, and S.-H. Myaeng. Identifying controversial issues and their sub-topics in news articles. In *Intelligence and Security Informatics*, pages 140–153. Springer, 2010.
- [7] M. Conover, J. Ratkiewicz, M. Francisco, B. Gonçalves, F. Menczer, and A. Flammini. Political Polarization on Twitter. In *ICWSM*, 2011.
- [8] D. L. Davies and D. W. Bouldin. A cluster separation measure. *IEEE TPAMI*, 1(2):224–227, 1979.
- [9] A. Doris-Down, H. Versee, and E. Gilbert. Political blend: an application designed to bring people together based on political differences. In C&T, pages 120–130, 2013.
- [10] K. M. Esterling, A. Fung, and T. Lee. How Much Disagreement is Good for Democratic Deliberation? The CaliforniaSpeaks Health Care Reform Experiment. SSRN, 2010.
- [11] S. R. Flaxman, S. Goel, and J. M. Rao. Filter bubbles, echo chambers, and online news consumption, 2015.
- [12] K. Garimella, G. De Francisci Morales, A. Gionis, and M. Mathioudakis. Exploring Controversy in Twitter. In *CSCW [demo]*, 2016.
- [13] E. Graells-Garrido, M. Lalmas, and D. Quercia. Data portraits: Connecting people of opposing views. arXiv preprint arXiv:1311.4658, 2013.
- [14] C. Grevet, L. G. Terveen, and E. Gilbert. Managing political differences in social media. In CSCW, pages 1400–1408, 2014.
- [15] P. H. C. Guerra, W. Meira Jr, C. Cardie, and R. Kleinberg. A measure of polarization on social media networks based on community boundaries. In *ICWSM*, 2013.
- [16] M. Jacomy, T. Venturini, S. Heymann, and M. Bastian. ForceAtlas2, a continuous graph layout algorithm for handy network visualization designed for the Gephi software, 2014.
- [17] G. Jeh and J. Widom. SimRank: A measure of structural-context similarity. In *KDD*, pages 538–543, 2002.
- [18] G. Karypis. CLUTO A clustering toolkit, 2002.
- [19] G. Karypis and V. Kumar. METIS Unstructured Graph Partitioning and Sparse Matrix Ordering System, 1995.
- [20] J. Kulshrestha, M. B. Zafar, L. E. Noboa, K. P. Gummadi, and S. Ghosh. Characterizing information diets of social media users. In *ICWSM*, 2015.
- [21] M. LaCour. A balanced news diet, not selective exposure: Evidence from a direct measure of media exposure. SSRN, 2012.
- [22] U. Maulik and S. Bandyopadhyay. Performance evaluation of some clustering algorithms and validity indices. *IEEE TPAMI*, 24(12):1650–1654, 2002.
- [23] Y. Mejova, A. X. Zhang, N. Diakopoulos, and C. Castillo. Controversy and sentiment in online news. Symposium on Computation + Journalism, 2014.
- [24] A. Morales, J. Borondo, J. Losada, and R. Benito. Measuring political polarization: Twitter shows the two sides of Venezuela. *Chaos*, 25(3), 2015.
- [25] S. A. Munson, S. Y. Lee, and P. Resnick. Encouraging reading of diverse political viewpoints with a browser widget. In *ICWSM*, 2013.
- [26] Y. Ruan, D. Fuhry, and S. Parthasarathy. Efficient community detection in large networks using content and links. In WWW, pages 1089–1098, 2013.
- [27] M. Thelwall. Heart and soul: Sentiment strength detection in the social web with sentistrength. In *CyberEmotions*, pages 1–14, 2013.
- [28] W. Zachary. An information flow model for conflict and fission in small groups. J. of Anthropological Research, 33: 452–473, 1977.