




Linking the Dynamics of User Stance to the Structure of Online Discussions

Christine LARGERON¹, Andrei Mardale¹, and Marian-Andrei RizoIU² 

¹ Univ Lyon, UJM -Saint -Etienne, CNRS, OGS, Saint-Etienne, France
Christine.LargerON@univ-st-etienne.fr

² Data Science Institute, University of Technology Sydney, Sydney, Australia
Marian-Andrei.RizoIU@uts.edu.au

Abstract. This paper studies the dynamics of opinion formation and polarization in social media. We investigate whether users' stance concerning contentious subjects is influenced by the online discussions they are exposed to and interactions with users supporting different stances. We set up a series of predictive exercises based on machine learning models. Users are described using several posting activities features capturing their overall activity levels, posting success, the reactions their posts attract from users of different stances, and the types of discussions in which they engage. Given the user description at present, the purpose is to predict their stance in the future. Using a dataset of Brexit discussions on the Reddit platform, we show that the activity features regularly outperform the textual baseline, confirming the link between exposure to discussion and opinion. We find that the most informative features relate to the stance composition of the discussion in which users prefer to engage.

Keywords: Online polarization dynamics · Online controversy · Social network analysis · Graph mining · Information diffusion

1 Introduction

In the twenty-first century, offline events are increasingly shaped by the discussions occurring on online social media. The outcome of significant events—such as the presidential elections in the United States of America [19, 20, 28] or the decision of the United Kingdom to leave the European Union [18]—were influenced by the opinions that voters formed using a wide array of online sources, including on social media.

Contentious subjects usually lead to heated arguments on social media, which in turn polarize public opinion. The prevailing theory is that online polarization emerges due to filter bubbles, which only expose users to peers with the same views [3]. This led to a body of work that believes that online polarization can be addressed by exposing users to contrary news and views [13, 16, 22]. However, participatory studies concluded that exposure to opposing views on social media could increase polarization [2, 21]. There is still an open gap concerning opinion formation on social media.

This work addresses two specific open questions concerning the dynamics of polarized opinion formation in the context of Reddit discussions around Brexit. The first open question deals with how users form polarized opinions. Some works claim that social media increases polarization [8,23], while other studies find that the usage of social media reduces polarization [4]. Furthermore, participatory and measurement studies [9] challenge this idea altogether, indicating that information savvy people leverage diverse sources of information and escape the filter bubble. The question is **are the stances of users concerning contentious subjects influenced by the discussions they are exposed to?** The second open question focuses on the dynamics of polarization. Previous work concentrates mainly on detecting and forecasting opinion polarization based on content diffusion in online social networks [11,27]; little work concentrates on detecting polarization dynamics. **Can we predict the future stance of users based on their present activity? and what are the factors that influence the changes of stances?**

We answer the questions mentioned above on a longitudinal dataset containing discussions around Brexit on Reddit, spanning from November 2015 until April 2019. Our work assumes two factors that determine users' stance towards contentious subjects. First, user stance has inertia, i.e., the stance at a given time is dependent on their past stance. Second, user stance depends on the stance of other users with whom the said user interacts. Consequently, the interactions with users of known stances indicate the future user stance, even without observing the textual content of these interactions.

We first divide the dataset time extent into fourteen time-periods, based on the notable events in the real-life Brexit timeline, such as the referendum, the triggering of Article 50, or the EU rejecting the UK's white paper. We investigate whether users' stance concerning contentious subjects is influenced by the online discussions they are exposed to and interactions with users supporting different stances. As there are no annotations available, we transfer a textual classifier trained on Twitter data to classify user stances in Reddit. Next, we answer the first open question by building three feature sets to describe user activity during each period. The purpose of these features is to capture a user's interaction with the other users of known stance in the community. The constructed features include overall activity levels, posting success, the reactions their posts attract from users of different stances, and the types of discussions in which users engage. We answer the second open question by setting up a series of predictive exercises that forecast the user stance in the next period based on the user description in the current period. We show that the activity features regularly outperform the textual baseline, indicating that user opinions are influenced by the discussions they are exposed to. We find that the discussion's stance composition that users prefer to engage in is the most informative feature. Notably, the content posted by a user during a time period appears to be less informative about the next period's user stance.

The main contributions of this work are as follows:

- We propose three feature sets predictive of the user stance that leverage solely the structure of the discussion (i.e., not the textual content emitted by the user).
- We show that all three feature sets are more predictive of the future stance than a textual baseline trained on the content emitted in the present.
- We provide predictive evidence that user polarization dynamics are linked to the stance composition of the discussions that the users are exposed to.

2 Related Work

We structure the discussion of the related works into two categories: detecting and alleviating polarization, and opinion and polarization dynamics.

Detecting and Alleviating Online Polarization. Previous work concentrates on detecting and reducing online polarization. Detection methods usually start from the social graph of the users. If the graph is presented as a signed network—i.e., the nodes are users, and the edges between users of the same polarity have a positive sign, while the edges across two polarities have a negative sign—community detection uncovers polarized communities [5]. The idea is to search for two communities (subsets of the network vertices) where within communities there are mostly positive edges while across communities there are mostly negative edges. When the sign of edges is not available, Garimella et al. [13] propose to use the diffusion cascades that occur on top of the social graph to detect the communities of users that participate together in the same cascades. Finally, they create a controversy score for discussion topics based on how polarized apart the communities are. In this work, we use a supervised approach to detect user polarization based on their emitted text: we use a textual classifier trained on annotated Twitter data to label the stance of Reddit users.

When it comes to reducing online polarization, it is generally assumed that exposing users to opposite views reduces their polarization [15,22]. Garimella et al. [12] devised tools and algorithms to bridge the polarized echo chambers and reduce controversy. They represent online discussions on controversial issues with an endorsement graph, and they cast the problem as an edge-recommendation problem on this graph. Graells-Garrido et al. [16] study how to take advantage of partial homophily to suggest agreeable content to users authored by people with opposite views on sensitive issues, while Musco et al. [25] search for the structure of a social network that minimizes disagreement and controversy simultaneously. However, empirical studies appear to contradict the fundamental thesis that users exposed to contrary views temper their polarization. Bail et al. [2] performed a participatory study on Twitter users, where they paid users to follow bots emitting tweets of the opposing opinion. They found that most users reinforced their previously held opinions and that exposure to opposing views on social media can increase political polarization.

Opinion and Polarization Dynamics. The prior work most relevant to this paper concerns the political polarization around Brexit and the study of polarization dynamics. Grčar et al. [17] studied the relation between the Twitter mood and the referendum outcome and who were the most influential Twitter users in the Pro- and Against- Brexit camps. They constructed a stance classification model, and they predicted the stance of about one million UK-based Twitter users. They found that the top pro-Brexit communities are considerably more polarized than the contra-Brexit camp. Amador Diaz Lopez et al. [1] collected 23 million Tweets related to the EU referendum in the UK to predict the Brexit vote. They used user-generated hashtags to build training sets related to the Leave/Remain campaign, and they trained an SVM to classify tweets. The above work uses textual content to decide the stance of a user. In contrast, our work leverages the structure of the discussion in which users engage without observing the textual content. In our experiments in Sect. 5 we show that our methods consistently outperform content-based methods.

When modeling the dynamics of opinion polarization, Das et al. [7] start from the conformity theory – i.e., a user will adopt the majority of their neighbors’ opinion – and propose a biased voter model. They show preliminary theoretical and simulation results on the convergence and structure of opinions in the entire network. On the empirical side, longitudinal study of controversy on Twitter [14] did not find long-term trends. However, they find that for particular subjects, polarization increased. By comparison, our work deals with the short-term polarization dynamics: we are interested in how users update their polarity concerning controversial topics based on their exposure to the content of different polarities.

3 The Dynamics of User Stance and Dataset

This section introduces our hypotheses around the dynamics of user stances, the structure of online discussions (on Reddit), and the dataset that we collected for the Brexit case study.

User Stance Dynamics. When faced with contentious subjects, users usually have opinions—dubbed here as *stances*; in the case of our study (i.e., Brexit) we define the following set of stances: **Against-Brexit**, **Pro-Brexit**, or **Neutral**. Our work’s central hypothesis is that users can update their stance as time passes by (for example, from Neutral to pro- or against- Brexit). Furthermore, we hypothesize that the change occurs partly due to the discussions the users are exposed to at present. We posit that stance changes occur on a much longer time scale than that of diffusions and threads. Without loss of generality, we assume that the time extent \mathcal{O} of our dataset is divided into periods (or time frames) during which each user’s stance is constant. The time periods are defined by a set of cutoff times $o_j \in \mathcal{O}$ such that $[o_j, o_{j+1}]$ defines an interval t . We denote the stance of a user u in a given time interval t as $c_t(u) \in \{A, P, N\}$. When passing from one period to the next, the users update their stance or maintain the stance from the previous period – in other words, the user u updates their stance from $c_t(u)$ to $c_{t+1}(u)$.

Structure of Discussion on Online Social Media. Online social networks can be viewed as meeting places where users have online discussions, submit content and articles in the form of text, link or media. In these meeting places, users interact with their peers, form and update opinions and stances towards topics. For example, on Reddit, users can start threads similar to forum environments or post comments on existing threads. Consequently, the discussions present themselves as hierarchies of posts in a tree-like structure. Figure 1a shows an example of a real Reddit discussion, containing an initial post (n_0) and five comments (n_1 to n_5). For instance, comment n_3 is a reply to comment n_1 . The resulted tree structure is shown in Fig. 1b, and leveraged in Sect. 4 to construct the non-textual features describing the activity of users. In the rest of this paper, we denote a tree of posts as *a thread*, which is started by *a post* – also known as *submission* in Reddit terminology, the root of the tree. We denote all the other nodes as *comments* – chronologically subsequent messages posted as replies to a post or other comments in the same thread.

We collectively denote posts and comments as *entries*. An entry is a triplet $s_j = (u_j, pc_j, d_j)$, where u_j denotes the user name, pc_j the published content, and d_j is the

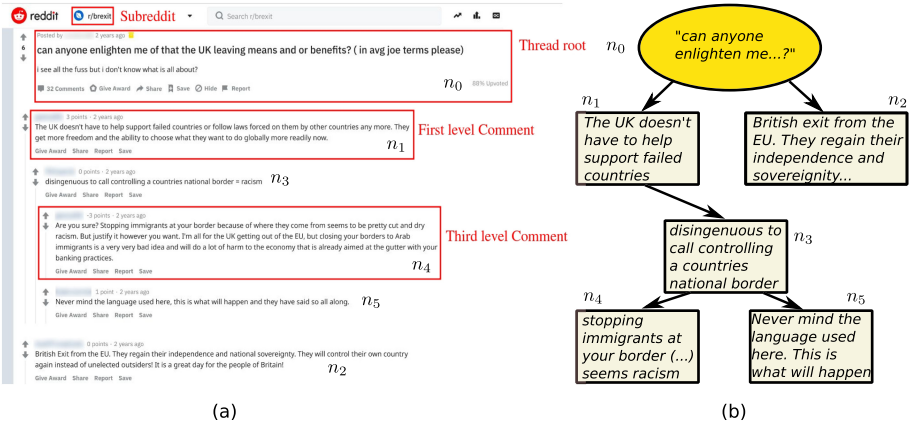


Fig. 1. (a) Elements of the Reddit platform. Structure of a discussion thread, with multi-level comments, inside a subreddit. (b) Logical structure used for analyzing the data.

time stamp of the entry s_j . We further define a *diffusion* δ_i as a temporally ordered sequence of entries, starting with a post, ending with a leaf comment and containing all comments on the path connecting the post and the leaf comment. Formally, it is defined as $\delta_i = \{s_j = (u_j, pc_j, d_j) | j = 1, 2, ..\}$. Visibly, there are as many diffusions in a thread as there are leaf nodes. For example, in Fig. 1b there are three diffusions: $\{n_0, n_1, n_3, n_4\}$, $\{n_0, n_1, n_3, n_5\}$ and $\{n_0, n_2\}$. Finally, a thread is a set of diffusions $S = \{\delta_i | i = 1, .., N\}$.

Dataset: Brexit Discussions on Reddit. We collected the Reddit dataset used in our case study using the Pushshift API [26]. It contains 229,619 entries (21,725 posts and 207,894 comments) posted between November 2015 and April 2019 on the *brexit* subreddit (<https://www.reddit.com/r/brexit/>). Each entry has the following variables: entry id, text, timestamp, author, parent id (useful for building the tree structure as shown in Fig. 1a), Reddit score, and the number of comments for the entry. A total of 14,362 unique authors participated in these discussions. We have divided the dataset’s time extent into 15 intervals based on the occurrence date of real events, such as the UK referendum of 23 June 2016, the nomination of M. Barnier as Chief Negotiator, beginning of the Brexit negotiations, rejection of the UK white paper by EU, the publication of the Brexit withdrawal agreement, first and second meaningful votes, etc. We split the entries into 15 subsets according to the time interval in which they were posted. Due to space constraints, we further profile the dataset and the 15 intervals in the online supplement¹. Also note that the constructed dataset, together with the code to build the feature sets detailed in Sect. 4.2 are publicly available².

¹ Supplementary Information available online: <https://arxiv.org/pdf/2101.09852.pdf#page=13>.

² Code and data publicly available: <https://github.com/behavioral-ds/online-opinion-dynamics>.

4 Forecast User Stance Dynamics

This section tackles the two research questions by posing them as supervised machine learning problems. We first describe the learning problem (Sect. 4.1); next, we build predictive features that embed user interactions with users of different stances (Sect. 4.2); finally, we describe the predictive setup (Sect. 4.4).

Table 1. Constructed feature sets describing user interactions with information diffusions.

| Feature set | Features |
|---------------------------------------|-----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| FS1 (User activity) | <ul style="list-style-type: none"> – number of initiated diffusions $ID_t(u)$ – number of submitted comments $CS_t(u)$ – quantiles of the number of received comments per entry $R_t^1(u), \dots, R_t^5(u)$ – stance at current time-frame $c_t(u)$ |
| FS2 (User activity per stance) | <ul style="list-style-type: none"> – number of comments submitted to entries from each stance $CS_t^A(u), CS_t^P(u), CS_t^N(u)$ – quantiles of the number of received comments per entry, tallied by commentator stance $R_t^{x1}(u), \dots, R_t^{x5}(u)$, $x \in \{A, P, N\}$ – stance at current time-frame $c_t(u)$ |
| FS3 (Structure of diffusion) | <ul style="list-style-type: none"> – quantiles of the number of submitted comments in diffusions per stance $UP_t^{x1}(u), \dots, UP_t^{x5}(u)$, $x \in \{A, P, N\}$ – stance at current time-frame $c_t(u)$ |
| FS4 (Relational features) | FS1 + FS2 + FS3 |
| FS0 (Textual features) | 100 top words + $c_t(u)$ |
| FS5 (All features) | FS0 + FS4 |

4.1 A Supervised Machine Learning Problem

We cast the problem of forecasting the future stance of users as a supervised machine learning problem. Each user u is represented by a set of features $FS_t(u)$ describing her Reddit activity during the time interval t . The feature set also includes the user stance at the current time t , i.e., $c_t(u)$. The task consists in forecasting the user stance at the next time interval $t + 1$, i.e., $c_{t+1}(u)$, using the features at time t , i.e., $FS_t(u)$. Off-the-shelf classifiers are used to learn a model from $FS_t(u)$ to $c_{t+1}(u)$. The difficulty lies in defining the features that describe the user’s activity during a period and obtaining the ground truth labels to build the training set and the test set. To determine user stances, we use a textual classifier trained on Twitter data (further detailed in Sect. 4.3). In the next section, we design several feature sets that capture users’ activity and their interactions with other users of different stances.

4.2 Predictive Features

We introduce three sets of features (denoted as **FS1**, ..., **FS4**, shown in Table 1) aimed at capturing increasingly complex information concerning user activity. **FS1** serves as an activity baseline, tallying user posting activity and the comments they receive. **FS2** aims to capture how the user interacts with users of different stances (e.g. do they prefer to comment on entries with similar stances to their own? to the opposite stance?), and whether they elicit more comments from users with the same polarity or the opposite. **FS3** aims to capture the type of threads in which the user engages (e.g., do they like to engage in discussion with a single stance or deliberative threads?). We detail each set.

FS1 focuses on the activity of the user at the global level. For a given user u and a time interval t , we count $ID_t(u)$, the number of diffusions initiated by u during the interval t (i.e., the number of posts sent by u) and $CS_t(u)$, the number of comments submitted by u during the interval t by excluding auto-comments. Thus, the number of entries submitted by u during the period is denoted $N_t(u)$ with $N_t(u) = ID_t(u) + CS_t(u)$. We also consider the user's success defined as the number of replies generated by his activity and quantified by the direct or indirect comments received by each entry (post or comment) submitted by u during the period. Formally, if r_i denotes the number of replies following the entry m_i submitted by u during the period t , we obtain the set $\{r_i | i = 1, \dots, N_t\}$ and we compute the quantiles $R_t^1(u), \dots, R_t^5(u)$ corresponding respectively to 0%, 25%, 50%, 75% and 100% of its distribution. Thus, **FS1** contains 8 features including $c_t(u)$ (the user stance at the current time).

FS2 aims to capture how the user interacts with users of different stances: **Against**, **Pro** or **Neutral** in our case study. First, we measure how the user engages with content from other users by counting the comments sent by user u during the period t in response to entries posted by each group denoted respectively $CS_t^A(u)$, $CS_t^P(u)$ and $CS_t^N(u)$. Thus, $CS_t(u) = CS_t^A(u) + CS_t^P(u) + CS_t^N(u)$, where $CS_t(u)$ has been defined in **FS1**. The underlying idea is to capture whether u exchanges more with users having the same stance as him or with users having a different stance. Second, we measure how the users of the different stances engage with u by counting the number of comments received from each group in response to entries sent by u during the period t . Thus, if r_i^x denotes the number of replies from group $x \in \{A, P, N\}$ following the entry m_i submitted by u during the period t , we obtain the distribution $\{r_i^x, i = 1, \dots, N_t\}$ and we compute the quantiles $R_t^{x1}(u), \dots, R_t^{x5}(u)$ corresponding respectively to 0%, 25%, 50%, 75% and 100% of this distribution. With this second set composed of 19 features, the objective is to capture whether content emitted by u attracts comments from the group of users of similar stance or from the other stances.

FS3 aims to capture the type of threads in which the user u engages. For each threads in which u posted an entry (post or comment) during the period, we compute the number of entries per group. More precisely, if NS denotes the number of threads in which u sent at least one entry during the period and S_i is one of these threads, we compute the number of entries A_i, P_i, N_i respectively emitted by each group in S_i . By this way, we obtain three sets $\{A_i | i = 1, \dots, NS\}$, $\{P_i | i = 1, \dots, NS\}$, $\{N_i | i = 1, \dots, NS\}$ that we can summarize by their respective quantiles $UP_t^{x1}(u), \dots, UP_t^{x5}(u)$, $x \in \{A, P, N\}$. Thus, if a user with a given stance, for example **Anti-Brexit**, prefers to exchange with the other anti-Brexit users, the features $UP_t^{A1}(u), \dots, UP_t^{A5}(u)$ will

Table 2. Hashtags used by Amador Diaz Lopez et al. [1] for splitting Twitter users in two categories, to train the Naive Bayes Classifier.

| Stance | Hashtags |
|-----------------------|------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| <i>Pro Brexit</i> | #voteleave #inorout #voteout #takecontrol #borisjohnson #lexit #independenceday #ivotedleave #projectfear #britain #boris #go #projecthope #takebackcontrol #labourleave #no2eu #betteroffout #june23 #democracy |
| <i>Against Brexit</i> | #strongerin #intogether #infor #votein #libdems #voting #incrowd #bremain #greenerin |

have higher values than $UP_t^{P1}(u), \dots, UP_t^{P5}(u), UP_t^{N1}(u), \dots, UP_t^{N5}(u)$. So, **FS3** contains 16 features, including $c_t(u)$. We also build **FS4**, the union of the above mentioned feature sets: $\mathbf{FS4} = \mathbf{FS1} \cup \mathbf{FS2} \cup \mathbf{FS3}$.

We also build a textual baseline based on the user’s content in the current period. We first extract the top 100 most frequent words (stop words removed) from all the Reddit dataset entries over all the time intervals. Next, we aggregate the text of all the entries of each user into a single document, and we compute the TF-IDF scores for the selected top 100 most frequent words. Consequently, **FS0** contains 101 features, including the user stance at the current time-frame $c_t(u)$. Finally, we also consider **FS5** composed of all the textual and relational features: $\mathbf{FS5} = \mathbf{FS4} \cup \mathbf{FS0}$.

4.3 Learning Stance in Twitter

One of the main challenges of this work is the lack of ground truth, *i.e.*, the stance for Reddit users at each time interval. We transfer to our Reddit dataset a model trained on Twitter and initially introduced by Amador Diaz Lopez et al. [1].

The Twitter Dataset. Amador Diaz Lopez et al. [1] collected the Twitter dataset from 6 January 2016 to July 2016 using the Twitter Firehose API. They crawled all the tweets using three search criteria related to Brexit: the general search term *Brexit*, hashtags such as *#leaveeu* or *#yes2eu* and Twitter usernames of groups and users set up to communicate about Brexit (e.g., *@voteleave* or *@yesforeurope*). The resulted dataset contains 26.5 million tweets emitted by 1.5 million users.

Build a Stance Predictor in Twitter. We build a Twitter stance predictor following the methodology proposed by Amador Diaz Lopez et al. [1], which we briefly summarize below. Amador Diaz Lopez et al. [1] curated two sets of hashtags, shown in Table 2, which indicate the user stance and utilized in 136 thousand tweets. We filter out occasional users – who emit less than 50 tweets – and users who do not employ any of the hashtags. For each of the remaining 11,277 users, we compute a ‘leave’ score equal to the difference between the number of used Pro-Brexit hashtags and Against-Brexit hashtags. We rank the users based on the score, and we select the 10% users with the lowest (negative) score as Against-Brexit users and the 10% of users with the highest (positive) score as Pro-Brexit users. The resulting set contains the aggregated tweets (one document per user) for 2,257 users. We first perform the usual text preprocessing: we remove stopwords, punctuation signs, hashtags, mentions, and other diacritics; we

convert all letters to lower case, remove rare words, and perform stemming. Next, we train a Naive Bayes classifier using 80% of the data, and we evaluate using the remaining 20% of the data. The model outputs the probability for a document (*i.e.*, a user) to belong to one of the classes (Against-Brexit or Pro-Brexit). Following the methodology proposed by Amador Diaz Lopez et al. [1], we convert this output probability into a discrete label: if the *leave* probability is below 0.25, we label the user as Against-Brexit; if it is greater than 0.75, we label the users as Pro-Brexit. Otherwise, the label is Neutral. On the test set, the trained model obtains a prediction macro-accuracy of 89.36% and a macro-F1-score of 88.68%. As shown in the next section, we transfer the trained model to compute $c_t(u)$, the users' stance in each period in the Reddit dataset. We use a Naive Bayes classifier because it is somewhat robust to concept drift and noisy features [29] – here, vocabulary change between Twitter and Reddit. The robustness is because rank scores are typically correct even if the conditional independence assumption is violated. We use cut-offs on the Naive Bayes score rather than interpreting the score as a probability in absolute terms.

4.4 Predictive Setup

Building Reddit Learning and Testing Sets. For each timeframe, we first aggregate all the Reddit messages of each user into a single document. Next, we assign them a Brexit stance using the Naive Bayes classifier trained on the Twitter dataset (detailed in Sect. 4.3). As we perform this procedure for each interval, we obtain not only the present stance of the user $c_t(u)$ but also the stance at the next timeframe $c_{t+1}(u)$. Finally, we compute the predictive feature sets **FS1**, ...**FS5** for each user and each period from the Reddit dataset.

Models and Evaluation. We predict users' stance in the next timeframe using each feature set computed on the current timeframe. We train and test five different algorithms – Logistic Regression, KNN, Random Forest, Gradient Boosting [10], XGBoost [6]. We evaluate using a double Cross-Validation. First, we use a 10-fold outer Cross-Validation to split the data into training and testing sets. At each fold, we tune hyper-parameters using an inner 5-fold Cross-Validation together with Random Search with 500 iterations. We measure performance using standard evaluation metrics and their standard deviation: macro-F1, macro-Accuracy, macro-Precision, and macro-Recall.

5 Results

This section presents the obtained performances for predicting the future stance of users. The Reddit dataset is imbalanced, with most of the user having a Neutral stance. Therefore, Fig. 2 plots the macro versions of accuracy and F1 score (macro-precision and macro-recall are shown in the Supplementary Information¹). Note that we use the macro version of the metrics, which gives equal representation to minority classes and alleviate the class imbalance in our dataset. Note that for a three-class classification problem (here Against, Neutral, Brexit), an unweighted, random classifier is expected to obtain an F1 score and accuracy score of 33%.

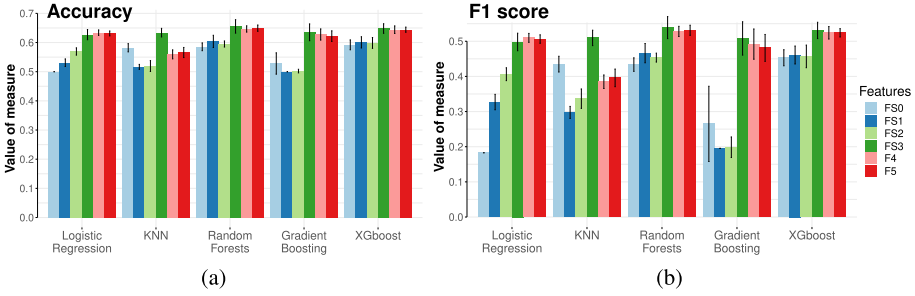


Fig. 2. Evaluation metrics for the developed models: accuracy (a) and F1-Score (b).

Table 3. F1 score of predicting next stance, tabulated per current stance.

| Stance at t | Stance at $t+1$ | | |
|---------------|-----------------|---------|--------|
| | Against | Neutral | Brexit |
| Against | 0.68 | 0.34 | 0.35 |
| Neutral | 0.51 | 0.62 | 0.45 |
| Brexit | 0.44 | 0.34 | 0.59 |

Analysis of the Relational Features. Figure 2 shows that the best classifier reaches 53.9% F1 score, which is double the random score. As the data is imbalanced, the accuracy is higher at 65%. For most classifiers, the performance improves from FS1 to FS3, with FS3 providing the best performance for all methods, except KNN. This indicates that the stance composition of the threads that the user prefers to engage in best indicates her future stance. The best performing classifiers are Random Forest and XGBoost. Interestingly, the combination of all activity features (denoted as FS4) does not further improve results.

Relational and Textual Features. Figure 2 shows that relational features (FS1 to FS4) have higher predictive power than textual features (FS0), for the best performing method (Random Forest and XGBoost). While the conclusions are more nuanced for the other classifiers, FS3 outperforms FS0 for all classifiers and all metrics. This result suggests that the type of discussions users engage in indicates their future stance more than the content they emit at present. Moreover, we observe that using textual together with relational features (FS5) does not improve results significantly as the performances of FS5 are equal to FS3.

Analysis per Stance. We analyze in more detail the performances of the best performing classifier (XGBoost) on the best features set (FS3). We compute the prediction performances for each combination of present and future stance – i.e., the nine combination $\{(c_t(u), c_{t+1}(u)) | c_t(u), c_{t+1}(u) \in \{A, B, N\}\}$. The values are reported in Table 3. We see that the classifier performs well for the users who maintain their opinion between two subsequent timeframes (shown by the main diagonal of Table 3). Noteworthy, it also performs well for the transitions from Neutral to Pro- or Against-Brexit, with F1

scores equal to 0.51 and 0.45, respectively. The result implies that we can predict the future stance of the currently undecided participants in online debates. The implications are significant, as most democratic processes tend to be decided by swaying undecided voters.

6 Conclusion

In this paper, we analyzed information diffusion in social media platforms, and we studied whether the stances of users are influenced by the discussions to which they are exposed. To capture the dynamics of the opinions of online communities, we chose the Reddit platform and Brexit as a case study due to its polarity. To better understand why users change their stance, we predict the future user stance using supervised machine learning algorithms. We construct three feature sets that capture different aspects of the user activity in the diffusion process. Our experiments showed that the best-performing feature set accounts for the stance composition of the threads in which a user chooses to engage. Notably, our activity feature sets outperform a textual baseline that encodes the content that the user emits.

One difficulty we met is the lack of ground truth, i.e., the stance for Reddit users at each time interval. To obtain the ground truth, we transferred a model trained on a Twitter dataset. However, the underlying distribution of language and structure of the two platforms differ. The transfer labeling risks introducing inaccuracies, and the performances would probably be better if the Reddit users' correct labels were available. This is a perspective of this work.

Acknowledgement. This work was partially supported by IDEXLYON ACADEMICS Project ANR-16-IDEX-0005 of the French National Research Agency, Facebook Research under the Content Policy Research Initiative grants, and the Defence Science and Technology Group of the Australian Department of Defence. We thank Keneth Benoit, who generously shared the Twitter dataset of Brexit discussions [1].

References

1. Amador Diaz Lopez, J.C., Collignon-Delmar, S., Benoit, K., Matsuo, A.: Predicting the Brexit vote by tracking and classifying public opinion using Twitter data. *Stat. Polit. Policy* **8**(1), 85–104 (2017). ISSN 2194–6299
2. Bail, C.A., et al.: Exposure to opposing views on social media can increase political polarization. *PNAS* **115**(37), 9216–9221 (2018)
3. Banisch, S., Olbrich, E.: Opinion polarization by learning from social feedback. *J. Math. Sociol.* **43**(2), 76–103 (2019). ISSN 0022–250X
4. Barberá, P.: How social media reduces mass political polarization. Evidence from Germany, Spain, and the US. *Midwest Pol. Sci. Assoc.* p. 44 (2014)
5. Bonchi, F., Gionis, A., Ordozgoiti, B., Galimberti, E., Ruffo, G.: Discovering polarized communities in signed networks. In: *CIKM*, pp. 961–970 (2019)
6. Chen, T., Guestrin, C.: Xgboost: a scalable tree boosting system. In: *KDD*, pp. 785–794. *ACM* (2016)
7. Das, A., Gollapudi, S., Munagala, K.: Modeling opinion dynamics in social networks. In: *WSDM*, pp. 403–412 (2014)

8. De-Wit, L., Brick, C., Van Der Linden, S.: Are social media driving political polarization? *Battles* (2019)
9. Dubois, E., Blank, G.: The echo chamber is overstated: the moderating effect of political interest and diverse media. *Inf. Comm. Soc.* **21**(5), 729–745 (2018)
10. Friedman, J.H.: Greedy function approximation: a gradient boosting machine. *Ann. Stat.* 1189–1232 (2001)
11. Garimella, K., De Francisci Morales, G., Gionis, A., Mathioudakis, M.: Quantifying controversy in social media. In: *WSDM*, vol. 1, pp. 33–42 (2016)
12. Garimella, K., De Francisci Morales, G., Gionis, A., Mathioudakis, M.: Reducing controversy by connecting opposing views. In: *WSDM*, pp. 81–90 (2017)
13. Garimella, K., Morales, G.D.F., Gionis, A., Mathioudakis, M.: Exposing twitter users to contrarian news (2017)
14. Garimella, K., Morales, G.D.F., Gionis, A., Mathioudakis, M.: The EBB and flow of controversial debates on social media. In: *ICWSM*, pp. 524–527 (2017)
15. Gillani, N., Yuan, A., Saveski, M., Vosoughi, S., Roy, D.: Me, my echo chamber, and I. In: *WWW*, pp. 823–831. *ACM* (2018)
16. Graells-Garrido, E., Lalmas, M., Quercia, D.: Data portraits: connecting people of opposing views. In: *International Conference on Intelligent User Interfaces* (2016)
17. Grčar, M., Cherepnalkoski, D., Mozetič, I., Kralj Novak, P.: Stance and influence of twitter users regarding the Brexit referendum. *Comp. Soc. Net.* **4**(1), 1–25 (2017)
18. Howard, P.N., Kollanyi, B.: Bots, #StrongerIn, and #Brexit: computationalpropaganda during the UK-EU referendum. Available at SSRN 2798311 (2016)
19. Hughes, A.L., Palen, L.: Twitter adoption and use in mass convergence and emergency events. *Int. J. Emerg. Manage.* **6**(3–4), 248–260 (2009)
20. Kim, D., Graham, T., Wan, Z., Rizoio, M.A.: Analysing user identity via time-sensitive semantic edit distance (t-SED): a case study of Russian trolls on Twitter. *J. Comput. Soc. Sci.* **2**(2), 331–351 (2019)
21. Liao, Q.V., Fu, W.T.: Beyond the filter bubble: interactive effects of perceived threat and topic involvement on selective exposure to information. In: *Human Factors in Computing Systems*, pp. 2359–2368 (2013)
22. Matakos, A., Terzi, E., Tsaparas, P.: Measuring and moderating opinion polarization in social networks. *Data Min. Knowl. Disc.* **31**(5), 1480–1505 (2017). <https://doi.org/10.1007/s10618-017-0527-9>
23. Messing, S., Westwood, S.J.: Selective exposure in the age of social media. *Commun. Res.* **41**(8), 1042–1063 (2014). ISSN 0093–6502
24. Mishra, S., Rizoio, M.A., Xie, L.: Modeling popularity in asynchronous social media streams with recurrent neural networks. In: *ICWSM*, pp. 1–10 (2018)
25. Musco, C., Musco, C., Tsourakakis, C.E.: Minimizing polarization and disagreement in social networks. In: *WWW*, pp. 369–378 (2018)
26. Pushshift: Pushshift. <https://pushshift.io/> (2019)
27. Rama, V., Garimella, K., Weber, I.: A long-term analysis of polarization on Twitter. In: *ICWSM*, pp. 528–531 (2017)
28. Rizoio, M.A., Graham, T., Zhang, R., Zhang, Y., Ackland, R., Xie, L.: #DebateNight: The role and influence of socialbots on twitter during the 1st 2016 us presidential debate. In: *ICWSM*, pp. 1–10 (2018)
29. Schütze, H., Manning, C.D., Raghavan, P.: *Introduction to Information Retrieval*. Vol. 39. Cambridge University Press Cambridge (2008)