

It's Always April Fools' Day! On the Difficulty of Social Network Misinformation Classification via Propagation Features

Mauro Conti, Daniele Lain,
Riccardo Lazzeretti, Giulio Lovisotto
University of Padua

Walter Quattrociocchi
IMT School for Advanced Studies, Lucca

Abstract

Given the huge impact that Online Social Networks (OSN) had in the way people get informed and form their opinion, they became an attractive playground for malicious entities that want to spread misinformation, and leverage their effect. In fact, misinformation easily spreads on OSN and is a huge threat for modern society, possibly influencing also the outcome of elections, or even putting people's life at risk (e.g., spreading "anti-vaccines" misinformation). Therefore, it is of paramount importance for our society to have some sort of "validation" on information spreading through OSN. The need for a wide-scale validation would greatly benefit from automatic tools.

In this paper, we show that it is difficult to carry out an automatic classification of misinformation considering only structural properties of content propagation cascades. We focus on structural properties, because they would be inherently difficult to be manipulated, with the aim of circumventing classification systems. To support our claim, we carry out an extensive evaluation on Facebook posts belonging to conspiracy theories (as representative of misinformation), and scientific news (representative of fact-checked content). Our findings show that conspiracy content actually reverberates in a way which is hard to distinguish from the one scientific content does: for the classification mechanisms we investigated, classification F_1 -score never exceeds 0.65 during content propagation stages, and is still less than 0.7 even after propagation is complete.

1 Introduction

An increasing number of people get informed on Online Social Networks (OSN) (Newman, N. and Levy, D.A.L. and Nielsen, R.K. 2015). However, as OSN allow every user to post content, which propagates among users through viral processes, these platforms became attractive targets for misinformation creators. Moreover, the hyperconnected world and increasing complexity of reality create a scenario in which viral processes on OSN are driven by confirmation bias, eliciting the proliferation of unsubstantiated rumors and hoaxes all the way up to conspiracy theories (Bessi et al. 2015b; 2015a). News stories undergo the same popularity dynamics as other forms of online contents (such as selfies and cat pictures) (Newman, N. and Levy, D.A.L. and

Nielsen, R.K. 2015). It is not a surprise then that the Oxford Dictionary in 2016 elected *Post-truth* as word of the year. The definition reads:

"Relating to or denoting circumstances in which objective facts are less influential in shaping public opinion than appeals to emotion and personal belief".

Several studies pointed out the effects of social influence online (Centola 2010; Fowler and Christakis 2010; Quattrociocchi, Caldarelli, and Scala 2014; Salganik, Dodds, and Watts 2006). Results reported in (Kramer, Guillory, and Hancock 2014) indicate that emotions expressed by others on Facebook influence our own emotions, providing experimental evidence of massive-scale contagion via social networks. As a result of disintermediated access to information and of algorithms used in content promotion, communication has become increasingly personalized, both in the way messages are framed and how they are shared across social networks. Selective exposure has been shown to favor the emergence of echo-chambers — polarized groups of like-minded people where users reinforce their world view with information adhering to their system of beliefs (Del Vicario et al. 2016). Confirmation bias, indeed, plays a pivotal role in informational cascades (Quattrociocchi, Scala, and Sunstein 2016; Bessi et al. 2015b; Zollo et al. 2015a; Sunstein 2002). Recent works (Bessi et al. 2015a; Zollo et al. 2015a) showed that attempts to debunk false information are largely ineffective. In particular, discussion degenerates when the two polarized communities interact with one another (Zollo et al. 2015b). OSN users therefore tend to select information that are consistent with their beliefs (even if containing false claims) and propagate it to like-minded friends, and to ignore information dissenting with their beliefs. This confirms that misinformation is a huge threat for modern society. Not only it can put people's life at risk, as in the case of "anti-vaccines" misinformation, it is also starting to be used against political opponents, such as in the US, during the 2016 electoral campaign (Silverman et al. 2016) and during Election Day (Rogers, K. and Bromwich, J.E. 2016).

Rising attention to the spreading of fake news and unsubstantiated rumors led researchers to investigate many of their aspects, from the characterization of conversation threads (Backstrom et al. 2013), to the detection of bursty topics on microblogging platforms (Diao et al. 2012), to

the disclosure of the mechanisms behind information diffusion for different kinds of contents (Romero, Meeder, and Kleinberg 2011). Spreading of misinformation also motivated major corporations like Google and Facebook to provide solutions to the problem (First Draft News Coalition 2016). However, given the amount of content posted everyday on OSN (e.g., Facebook reports 1.18 billion daily active users on September 2016 (Facebook 2016)), effective fact-checking would greatly benefit of automatic classification tools, that would possibly not require human intervention. Moreover, classification of fake news and misinformation should ideally use properties that misinformation creators can not manipulate. Considering in the classification, for example, the trustworthiness of news domains, or the topic of content, could lead to an arms race with creators of false news. Instead, the topology of propagation cascades, and the patterns of users' interaction with content, are outside of the domain of our "adversaries", and are much more difficult to be manipulated.

In this work, we investigate detection of viral processes by comparing diffusion of posts from scientific and conspiracy pages on the Italian Facebook network. The former diffuse scientific knowledge, where details about the sources (such as authors and funding programs) are easy to access. We therefore select their posts as representative of fact-checked content. The latter aim at diffusing what is neglected by "manipulated" mainstream media. Specifically, conspiracy theories tend to reduce the complexity of reality by explaining significant social or political aspects as plots conceived by powerful individuals or organizations. Since these kinds of arguments can sometimes involve the rejection of science, alternative explanations are invoked to replace the scientific evidence. For instance, people who reject the link between HIV and AIDS generally believe that AIDS was created by the US Government to control the African American population. We therefore select their posts as representative of misinformation.

Contribution. The contributions of this paper are the following:

- We show that automatic fact-checking with classification techniques employing only structural features of content propagation cascades (features that are robust to attacker's manipulation) does not suggest to bring usable results. Given the grain of our data, we design classifiers that leverages topological properties of content propagation cascades, and properties of the evolution over time of users' interactions with content. A set of classifiers operates with evolution properties to classify content during its early propagation stage, another set of classifiers operates with more features to classify content that already propagated. Indeed, being able to issue warnings about possible fake news as early as possible, and retroactively flag such news can be useful, in the fight against OSN misinformation.
- We evaluate our classifiers on a well-known dataset of Facebook posts from Italian pages. We use posts belonging to conspiracy theories as representative of misinfor-

mation, and posts belonging to scientific news as representative of fact-checked content. Our findings highlight the complexity of creating automatic solutions to misinformation classification. Indeed, structural features of content propagation do not allow us to obtain notable improvements from a random guess baseline: F_1 -score for the classification mechanisms we investigated is always lower than 0.65 during content propagation stages, and is still less than 0.7 even after propagation is complete.

Outline. Section 2 overviews related work in news and misinformation propagation in OSN. Section 3 formally models and defines content propagation in Facebook. Section 4 presents our experiment design and the features we extract from content propagation cascades. Then, Section 5 evaluates our classification and, finally, Section 6 critically discusses our results, summarizes the paper and delineates future work.

2 Related Work

Several studies moved towards the spreading of rumors and behaviors on online social networks, challenging both their structural properties and their effects on social dynamics (Moreno, Nekovee, and Pacheco 2004; Doerr, Fouz, and Friedrich 2012; Seo, Mohapatra, and Abdelzaher 2012; Borge-Holthoefer et al. 2013; Cozzo et al. 2013; Borge-Holthoefer, Rivero, and Moreno 2012). In (Ugander et al. 2012), authors find that the probability of contagion is tightly controlled by the number of connected components in an individual's contacts neighborhood, rather than by the actual size of the neighborhood. In (Centola and Macy 2007) researchers show that, although long ties are relevant for spreading information about an innovation or social movement, they are not sufficient with respect to the social reinforcement necessary to act on that information. A key factor in identifying true contagion in social network is to distinguish between peer-to-peer influence and homophily: in the first case, a node influences or causes outcomes to its neighbors, whereas in the second one, dyadic similarities between nodes create correlated outcome patterns among neighbors that could mimic viral contagions even without direct causal influence (McPherson, Smith-Lovin, and Cook 2001). The study presented in (Aiello et al. 2012) reveals that there is a substantial level of topical similarity among users which are close to each other in the social network, suggesting that users with similar interests are more likely to be friends. In (Aral, Muchnik, and Sundararajan 2009) authors develop an estimation framework to distinguish influence and homophily effects in dynamic networks and find that homophily explains more than 50% of the perceived behavioral contagion. In (Bakshy et al. 2012) the analysis faces the role of social networks and exposure to friends' activities in information resharing on Facebook. Once having isolated contagion from other confounding effects such as homophily, authors claim that there is a considerably higher chance to share contents when users are exposed to friends' resharing. All these contributions strive to understand the inner mechanism of rumor spreading and to eventually predict

massive diffusion processes, i.e. cascades. Cascades recurrence and prediction has been shaped in (Cheng et al. 2014) and (Cheng et al. 2016).

3 Methods

In this section, we first report and describe the employed dataset (Section 3.1). We then give some necessary background knowledge, and present our reference model and definition of content propagation mechanisms of Facebook (Section 3.2).

3.1 Dataset

We employ a well-known dataset of posts shared by Italian Facebook users (Bessi et al. 2015a). This dataset contains posts published by 73 public Facebook pages: 34 pages that publish scientific content (e.g., press releases of peer-reviewed articles), and 39 pages that publish conspiracy theories-related content (e.g. new world order, chemtrails). Additionally, the dataset contains information about the interaction of users with these posts, and users’ ego-networks (i.e., the list of users that are their friends, when such list is public). Additionally, for a set of posts, the dataset provides information about the *propagation cascade* of such content, generated by users’ reshares, and subsequent reshares from their friends. This propagation from one user to other users can happen multiple times, forming a cascade of resharing. Using information from the dataset, we extracted 112141 non-empty propagation cascades, 89491 for conspiracy and 22650 for science, respectively. We underline that the dataset is obtained by using the Facebook Graph API, and contains only public information. Hence, timestamps of *reshares* and *comments* are available, but timestamps of *like* interactions are not.

3.2 Background and Definitions

We now present our reference formal representation of Facebook’s *friendship graph*, and the *potential propagation graph* generated by content posted on the social network.

Facebook Friendship Graph. We model Facebook relationships as a graph $\mathcal{G}\langle V, E \rangle$, that we call *Facebook friendship graph*, where V is the set of nodes that represent *entities*, namely user accounts and page accounts. We assume two main differences among these two types of entity: (i) pages can post new content on the OSN, while users can only interact with such content by liking, commenting, and resharing it; and (ii) users can establish *friendship* relationships with other users, while pages cannot. Indeed, two users $v_1, v_2 \in V$ are connected by an edge $e(v_1, v_2) \in E$ if they are *friends* on Facebook. Pages are not connected by any edge, as they do not have proper friendship relationships. This model is a simplification of how Facebook actually works, because users can post new content, and pages and users are linked by *like* relationships. Similarly, users can *follow* other users, without having any friend relationship with them. However, for this work, we do not focus on new content generated by users. Moreover, the dataset

we use lacks information about the *like* and *follow* relationships, that we therefore can not consider.

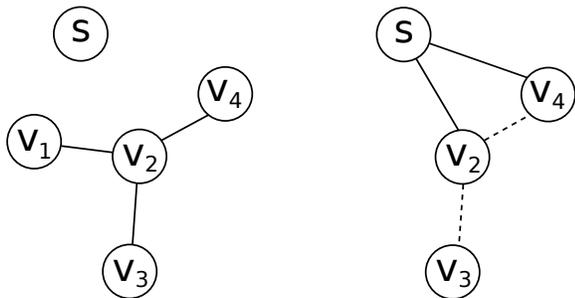
Potential Propagation Graph. Before formally modeling the spreading of content on Facebook, we describe some fundamental concepts of the OSN. We recall that, in our model, only pages can post new content on the social network. Henceforth, we refer to the page that originally posted some content as the *seed page*. Instead, users find new content posted by the pages they like, and content that their friends recently interacted with. They can then interact with such content through the means of *resharing*, *liking*, and *commenting* it. However, all of these interactions can happen directly on the original content, or on some types of interactions of users’ friends. For example, a user observing a comment or reshare of a post by one of his friends can decide to interact directly with the original post, or with his friend’s interaction itself. In the first case, user’s interaction looks exactly the same and it is impossible to understand whether the content was found thanks to his friend’s interaction, or directly on the seed page. In fact, differently from related work (Friggeri et al. 2014), the dataset we employ does not provide this type of information. We therefore take a conservative approach, saying that the content *potentially propagated* to the user both from the seed page and from any of his friends that interacted with the content, without any distinction. On the other hand, in the second case, it becomes clear that the user found the content thanks to his friend. We therefore say that the content *propagated* to the user from his friend.

We now formalize the above observations. Let P be the set of contents posted on Facebook by seed pages. For each post $p \in P$, created at some time t by a seed page, at a generic subsequent point in time $t + \delta$ we define a *potential propagation graph* $\mathcal{G}_{t+\delta}^p\langle V_{t+\delta}^p, E_{t+\delta}^p \rangle$, where $V_{t+\delta}^p$ is the set composed by the seed page, and by the users that interacted with p during the time interval $[t, t + \delta]$. *Final potential propagation graph* $\mathcal{G}^p\langle V^p, E^p \rangle$ is the graph formed considering all interactions with p on the timespan of the analysis (as new interactions with old content are always possible on Facebook). Two nodes $v_1, v_2 \in V_{t+\delta}^p$ are connected by an undirected *potential propagation edge* $e^p(v_1, v_2) \in E_{t+\delta}^p$ if either (i) v_1 or v_2 already interacted with p and $\exists e \mid e(v_1, v_2) \in E$ (that is, v_1 and v_2 are friends on Facebook), or (ii) v_1 or v_2 is the seed page. Therefore, an edge $e(v_1, v_2) \in E_{t+\delta}^p$ indicates that the content p *potentially propagated* either from v_1 to v_2 , or from v_2 to v_1 . We associate to each edge $e \in E_{t+\delta}^p$ two different properties:

1. **Time:** the time when the interaction between v_1 and v_2 took place;
2. **Type:** either “like”, “comment”, “reshare”, or “friendship”, depending on the type of interaction between v_1 and v_2 .

Hereafter in the paper, we will refer to these properties for an edge $e \in E_{t+\delta}^p$ with *e.property_name* (for example, *e.time*).

Figure 1 depicts an example of this propagation model. We represent a simple Facebook friendship graph in Figure 1a, where edges represent friendship relations. Nodes v_1, \dots, v_4 are OSN users, while node s represents a public page posting content on the platform. We suppose v_2 and v_4 interacted with a given content posted by s , and v_3 re-shared the post from v_2 . Node v_1 did not interact with the content, hence is not present on the potential propagation graph of the post, that we represent in Figure 1b. There, with edges $(s, v_2), (s, v_4), (v_2, v_4)$ we represent possible propagation paths of the content: either v_2 or v_4 could have seen the content from the seed page, or from previous interaction of their friend. Node v_3 only has an edge (v_2, v_3) , as we know that the interaction happened thanks to v_2 , and therefore the content propagated from this node. Additionally, we highlight that some possible propagation edges, represented with a dashed line, correspond to friendship edges of the friendship graph.



(a) Facebook Friendship Graph (b) Potential Propagation Graph

Figure 1: Sample Facebook friendship graph, and potential propagation graph of some content. A dashed line represents potential propagation edges that also corresponds to a friendship edge of the friendship graph. Node v_1 does not interact with the content, and is therefore not present in the potential propagation graph.

4 Experiments

In this section, we describe the details of the experiment and of the analysis performed on the dataset. We first describe the experimental design (Section 4.1). We then thoroughly report and motivate the different features that we extracted from propagation cascades (Section 4.2).

4.1 Experimental Design

The aim of our experiments is to show that it is difficult to discriminate between conspiracy theories, and fact-checked scientific news, by using only the propagation graph of the post, as it would be difficult to manipulate by misinformation creators. In particular, we evaluate this in two specific moments of content propagation:

1. **Early Stage**, meaning that classification of the type of post happens as early as possible during its propagation phase;
2. **Final Stage**, meaning that classification happens when the post already stopped propagating (within the considered timespan of the employed dataset), and its cascade is complete.

These two scenarios are of particular interest, both from a research and a practical points of view. Indeed, being able to issue automated warnings as early as possible about potential misinformation could help in reducing their spread (Friggeri et al. 2014). On the other hand, still being able to perform such a classification when the post already propagated might serve as a warning in order to prevent its further diffusion. Moreover, retroactively flagging old posts as potentially fake could help training users to discriminate between fact-checked and dubious information, a major direction in the fight against misinformation (First Draft News Coalition 2016).

To investigate these two scenarios, we set up two different experiments, modeled as binary classification tasks, where the two classes are *conspiracy* and *science*. We describe them in the following.

Early Stage Classification. In this scenario, we have access only to users’ interactions with a post up to a certain time $t + \delta$, after it has been created at time t , (e.g., up to after two hours after creation). With this scenario, we simulate how OSNs such as Facebook could try to continuously detect misinformation content during its propagation. Therefore, for each post p , given its creation time t , we extract features from the partial propagation graph $\mathcal{G}_{t+\delta}^p$, and from the evolution of properties of the propagation graph before the current time $t + \delta$ (discussed further in Section 4.2). We use 30 minutes steps, as this proved to be a good trade-off between the granularity of the analysis and the number of intervals to consider. We stop at two days (2880 minutes, or 96 time steps) after the publication time of a post, because we observed that most of the interactions happen in this period ($> 95\%$). Finally, we compare the performance of different well-known classifiers, namely Random Forests (Ho 1995), Linear Discriminant Analysis (Izenman 2013), and Multi-Layer Perceptron (Haykin 1998), for each $\delta \in [30, 60, \dots, 2880]$, in a cross-validation scheme.

Final Stage Classification. In this scenario, we suppose that a post already stopped its propagation, and no new interactions with it have been observed for a long time. We describe the final propagation graph \mathcal{G}^p of post p with a set of high-level and topological features, that we describe in more detail in Section 4.2. We then compare the performance of different classifiers, in a cross-validation scheme.

4.2 Feature Extraction

To extract information from the different propagation graphs \mathcal{G}^p and \mathcal{G}_t^p , we identify three possible categories of features: (i) high-level properties of the content propagation, (ii) topological properties of the propagation graphs, (iii) evolution

properties of the content propagation. We use feature sets (i), (ii), and (iii) considering the evolution after two days, in the Final Stage Classification scenario. We only use feature set (iii), in the Early Stage Classification scenario, because the other features characterize only \mathcal{G}^p , but not \mathcal{G}_t^p . We summarize these different features in Table 1, along with their formal definitions, and describe them in the following.

High Level Properties. These features represent high-level properties of the complete propagation cascade represented by $\mathcal{G}^p(V^p, E^p)$.

Some features represent very general properties related to the virality of the content. These general properties are the *lifetime* of the cascade, measured as the distance in minutes from the first to the last captured interaction with the content; the *size* of the cascade in terms of number of nodes (users who interacted with the content); the *number of total interactions*; and the time required for the cascade to reach its 90% total interactions (referred to as *90% interactions time*).

Other features derived from high level properties attempt to capture different possible types of interaction with the content. *Friendships ratio* is defined as the proportion of edges whose type is “friendship” over the total number of edges and represents the number of times, in proportion, that the post potentially propagated among friends, rather than directly from the seed node. Indeed, if no friends of users interact with some content, its potential propagation graph only contains edges with the seed page. *Interactions ratio*, instead, represents the average exposure to interactions from friends of users with the content. Since vertices are interacting users, and edges are potential interactions, higher values of this metric mean lower exposure (little interaction with the content by one’s friends). These features are motivated by the observation that it is possible that the users’ fruition of different types of content is different, with some types of content being interacted directly from the source, and other types of content relying more on word-of-mouth propagation (Romero, Meeder, and Kleinberg 2011).

Topological Properties. We also select as features some well-known properties of the topological structure of graphs. These properties are commonly used to learn information about graph structures, and have been applied in solving problems such as link analysis and prediction (Al Hasan et al. 2006), and especially in cascade and virality prediction (Hong, Dan, and Davison 2011; Cheng et al. 2014). The *average vertex degree* feature represents the average number of possible propagation edges for the content at any given hop. Higher values of this metric indicate the presence of interacting users greatly exposed to the content, or able to influence many of their social friends. The *global clustering coefficient* (Holland and Leinhardt 1971), a measure of the density of connections of graphs, is another indication of whether the possible propagation paths are generated from interactions between friends, or directly with the seed page. *Assortativity coefficient*, defined as the degree correlation between pairs of linked nodes (Newman 2002), can measure how friends influence each others in interacting with content on the social network. *Average path length* (Fronczak,

and Hołyst 2004), also known as Wiener index, gives us indications of the virality of the content, in terms of distance of propagation from the seed page. Long cascading news, reshared many times from interacting friends, will exhibit a longer average path length than news whose interactions happened mostly from the seed node. Finally, the *diameter* of a graph, defined as the longer shortest path between any pair of nodes of the graph, indicates the spreading distance of posts.

Evolution Properties. These features represent evolution properties over time of the post p propagation, from its creation time t to a subsequent point $t + \delta$. To compute these features, we construct the propagation graphs at different time steps $\mathcal{G}_{t+30}, \dots, \mathcal{G}_{t+\delta}$. We then calculate the value of three of our high-level features for each graph at each time step: (i) Friendships Ratio, (ii) Size, and (iii) Interactions Ratio. For each of these high-level feature, we obtain a time series

$$v_1, v_2, \dots, v_{\delta/30},$$

on which we compute a set of well-known statistical measures that represent the evolution over time of the series (Wiens, Horvitz, and Gutttag 2012). Namely, these statistical measures are the *mean*, *standard deviation*, *linear weighted mean* and *quadratic weighted mean*, *average absolute change*, and *maximum* of the series. These measures capture the evolution of the time series up to the current time and, especially, do not require us to know the whole time series, an important property for our Early Stage classification experiment. We decided to derive time-series only for the three high-level features listed above. Indeed, we argue that the other features either have no temporal properties (e.g., lifetime, time to reach 90% interactions), evolve in similar or predictable ways (e.g., diameter), or describe behaviors that are already captured by the selected time series. Moreover, evolution properties of time series would be especially used in Early Stage classification. Social networks need to perform feature extraction on this scenario at every time step: features derived from topological properties are too computationally expensive (e.g., diameter calculation runs in more than quadratic complexity w.r.t. the number of vertices) to be continuously calculated for each new post.

5 Results

In this section we present the results of our experiments. In particular, we first discuss evaluation metrics and baseline (Section 5.1). We then report the results on the Early Stage Classification scenario (Section 5.2), and the results on the Final Stage Classification scenario (Section 5.3).

5.1 Evaluation Metrics

As usual in binary classification, the classification baseline is the performance of a random classifier on the data: without any information regarding the propagation graph, it guesses either *science* or *conspiracy* with equal probability, as a coin flip. The goal of our experiments is to show that structural features do not help more sophisticated models in improving the baseline performance.

Feature Name	Description
High Level Properties	
Size	Number of edges of the graph $ E^p $
Friendships Ratio	Proportion of edges whose type is “friendship”: $ \{e \in E^p \mid e.type = \text{“friendship”}\} / E^p $
Interactions Ratio	Number of vertices over the number of edges: $ V^p / E^p $
Lifetime	Time passed between post creation time and the last interaction time: $\max_{e \in E^p} e.time - t$
90% Interactions Time	Time required for the content to reach its 90% number of interactions
Topological Properties	
Average Vertex Degree	Average possible propagation paths of the post: $2 \cdot E^p / V^p $
Clustering Coefficient	Ratio of connected triplets of nodes
Assortativity Coefficient	Degree correlation between pairs of linked nodes
Average Path Length	Average length of the shortest paths
Diameter	Length of the longest shortest path between any pair of vertices
Evolution Properties	
Mean	$\frac{1}{n} \sum_i^n v_i$
Linear Weighted Mean	$\frac{2}{n(n+1)} \sum_i^n i \cdot v_i$
Quadratic Weighted Mean	$\frac{6}{n(n+1)(2n+1)} \sum_i^n i^2 \cdot v_i$
Standard Deviation	$\sqrt{\frac{\sum_i^n (v_i - \bar{v})^2}{n}}$, where \bar{v} is the mean of v
Average Absolute Change	$\frac{1}{n} \sum_i^{n-1} v_i - v_{i+1} $
Maximum	$\max_i v_i$

Table 1: Description and formal definition of features we use on our Final Stage and Early Stage classification experiments.

Unfortunately, our dataset is highly imbalanced (composed by 89491 news for *conspiracy*, and only 22650 news for *science*). With such imbalance, standard evaluation metrics (such as *precision*, *accuracy*, and *recall*) can be misleading, because they do not account for the uneven class frequencies. Even computing averages of these metrics using weights based on the class frequencies does not fit our intentions of consistently comparing our results with the fixed baseline.

To deal with this imbalance, we performed two distinct experiments: (i) we consider only metrics that take imbalance into account and use the full dataset, and (ii) we consider meaningful metrics with balanced dataset, obtained *undersampling* the original dataset.

To perform (i), as metrics we use Area Under Receiver Operating Curve (AUC) (Hanley and McNeil 1982), and Cohen’s Kappa (Cohen 1968) (scaled into the interval $[0, 1]$). Indeed, the value of these metrics for a random classifier is exactly 0.5, which we use as baseline. In this way, we can use the full dataset and still be able to compare our results

with the baseline.

To perform (ii), we undersampled the most frequent class (i.e., *conspiracy*). We therefore extracted exactly 22650 *conspiracy* samples from the full dataset, and created a subsample with perfectly balanced classes. To account for possible biases caused by the undersampling, we repeated the process several times, and averaged the outcomes. Using perfectly balanced datasets, we can evaluate *precision*, *recall*, *accuracy*, and F_1 -score values of our classifiers. Indeed, the value of these metrics for a random classifier on balanced data is exactly 0.5, which we use as baseline.

Hereafter, if not differently specified, the metrics that require a positive class (such as *precision*, *recall*, and F_1 -score) use *conspiracy* as the positive class, and *science* as the negative one.

5.2 Early Stage Classification

To evaluate this scenario, we recall that we used a total of 18 features. We then followed the methodology explained in Section 4.1, and evaluated three separate clas-

sifiers (namely, Linear Discriminant (LD), Random Forest (RF), Multi-Layer Perceptron (MLP)) at every 30 minutes time step. As discussed, we first evaluated AUC and Cohen’s Kappa of these classifiers on the full, unbalanced dataset. We then evaluated precision, recall, and accuracy of these classifiers on undersampled balanced datasets.

Figure 2 reports the AUC and the Cohen’s Kappa metrics, on a 5-fold cross-validation scheme, at different time steps after the original post is published. We can see that evolution properties of the propagation graph do not help significantly our classifiers. Indeed, the curves are close to the random classification baseline (dotted horizontal line), and far from perfect classification (i.e., 1.0). Cohen’s Kappa lies almost exactly on the random classification baseline, suggesting that the classification performance is almost random. Using AUC, it is possible to set good classification thresholds, slightly improving the performance. MLP performs slightly better than RF and LD, and we can observe that results do not change significantly after a few hours of analysis. However the outcome remains unsatisfying (i.e., AUC under 0.6 for each classifier at each time step).

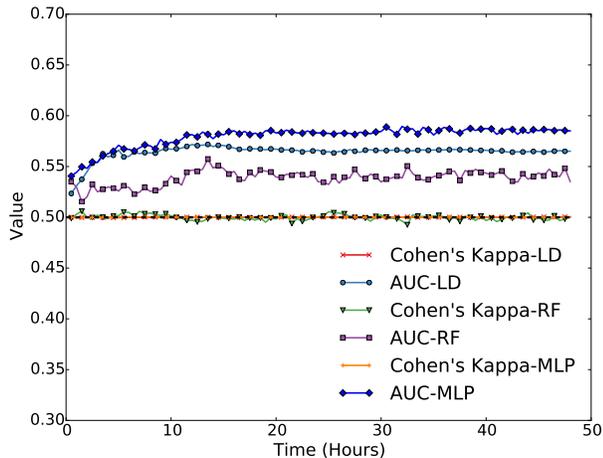


Figure 2: Early Stage Classification scenario, AUC and Cohen’s Kappa, as a function of elapsed time.

Figure 3 reports the F_1 -score obtained on the undersampled balanced dataset, on a 5-fold cross-validation scheme, averaged over ten repetitions of the undersampling procedure afterwards. Figure 3 further confirms the findings highlighted in Figure 2: classifiers are not able to leverage the evolution properties to discriminate between *science* and *conspiracy*. In fact the F_1 -score curve, after a slight increment during the first 10 hours, stabilizes relatively close to the baseline, and remains below 0.65 throughout the whole 48 hours timespan.

5.3 Final Stage Classification

To evaluate this scenario, we recall that we used a total of 28 features with three distinct classifiers: Linear Discriminant (LD), Random Forest (RF), Multi-Layer Perceptron (MLP). Again, we first evaluated the unbalanced dataset, and then evaluated the undersampled balanced datasets.

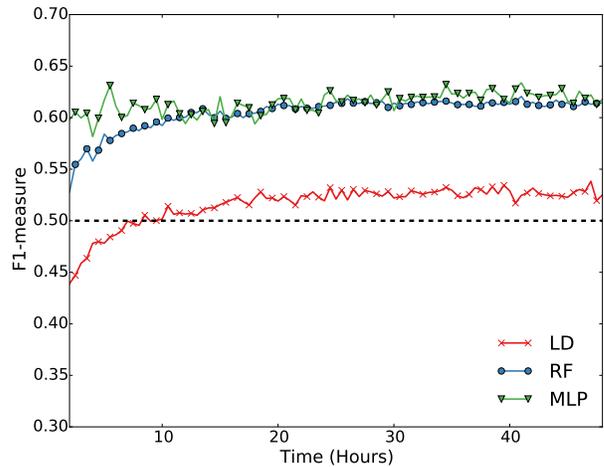


Figure 3: Early Stage Classification scenario, classification F_1 -score, as a function of elapsed time.

Figure 4 reports average AUC and Cohen’s Kappa on a 5-fold cross-validation scheme, on the full dataset. The dotted horizontal line shows the baseline performance of random classification. Indeed, we observe that our classifiers do not significantly improve the baseline, with the metrics remaining below 0.75. LD performance is poorer than RF and MLP, probably due to the simplicity of the classification tool.

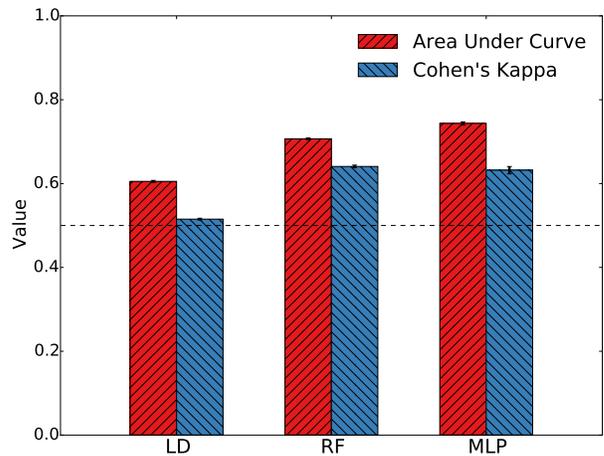


Figure 4: Final Stage Classification scenario, AUC and Cohen’s Kappa.

In Figure 5 we report the ROC curve, on a 5-fold cross-validation scheme, on the full dataset. From Figure 5 we can observe that no classifier can reach a good tradeoff between true positive rate and false positive rate. Indeed, the curves are relatively close to the baseline (diagonal dotted line), meaning that, as the decision threshold changes, lots of samples are misclassified.

Table 2 reports several metrics computed on the undersampled balanced dataset. Results are averaged on a 5-fold

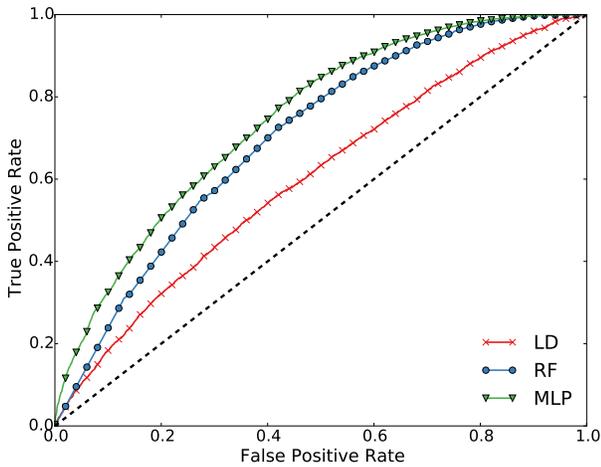


Figure 5: Final Stage Classification scenario, ROC curves. For visual clarity, we linearly interpolated the points in the curves, and plotted markers at each 0.02 step in the false positive rate.

cross-validation scheme, and then over ten repetitions of the undersampling procedure. These metrics show that even if the classifiers are able to improve the baseline slightly, they can not reach good performance.

Classifier	Precision	Recall	Accuracy	F_1 score
LD	0.578	0.523	0.570	0.549
RF	0.654	0.742	0.675	0.695
MLP	0.659	0.688	0.665	0.672

Table 2: Final Stage Classification scenario, performance of the classifiers.

6 Conclusions

Early detection of misinformation plays a crucial role in social networks. In this paper, we analyzed the difficulty of discerning conspiracy posts from scientific posts on Facebook. We focused on using only structural features of content propagation, because they cannot be easily manipulated by misinformation creators. Our results show that misinformation classification during its early propagation stage with these features is unsuccessful, suggesting that the spreading dynamics captured by our features are independent on the type of content. Furthermore, even considering the cascade at the end of content propagation does not help: also in this case, the improvement provided by a classifier over random coin flips is negligible.

Our findings suggest that in Facebook users interact with different types of content in similar ways, reinforcing the hypothesis of echo chambers (Del Vicario et al. 2016). Inside these chambers, strongly polarized by topic (Bessi et al. 2015b), content propagation exhibits very similar structural properties, that are therefore less useful in content clas-

sification. These results highlight the necessity of including content-related features, or polarization metrics, in future analysis (i.e., whether particular users and their echo chambers are more polarized towards one type of content). Unfortunately, misinformation creators can easily control content-related features, in order to avoid algorithmic detection. Moreover, user polarization can be clearly understood from past users’ behaviors, but it takes time to understand polarization of new users. Hence, automatic detection of fake news remains an open challenge.

Future Work. The employed dataset has some limitations, bound to the Facebook API: (i) it only contains Facebook public information; (ii) it does not contain the timestamp of likes, one of the most common interactions; and (iii) it does not always provide information about whether interaction with content happened because of interactions of user’s friends. We would like to analyze finer-grained data, that takes these factors into account, because it could lead to improved results.

In our analysis, we identified a set of features and used them in well-known classifiers. Our experiments were extensive, but not complete. However, we expect that the use of different models would not provide significant improvement. This claim needs to be further validated, possibly also using more recent datasets.

In the future, we also plan to test different methods for misinformation classification, based on user polarization and content-related features, to investigate whether these information could help propagation properties, and overcome the difficulty of this problem.

References

- Aiello, L.; Barrat, A.; Schifanella, R.; Cattuto, C.; Markines, B.; and Menczer, F. 2012. Friendship prediction and homophily in social media. *ACM Transactions on the Web (TWEB)* 6(2):1–33.
- Al Hasan, M.; Chaoji, V.; Salem, S.; and Zaki, M. 2006. Link prediction using supervised learning. In *International Conference on Data Mining (SDM): Workshop on Link Analysis, Counter-Terrorism and Security*.
- Aral, S.; Muchnik, L.; and Sundararajan, A. 2009. Distinguishing influence-based contagion from homophily-driven diffusion in dynamic networks. *Proceedings of the National Academy of Sciences (PNAS)* 106(51):21544–21549.
- Backstrom, L.; Kleinberg, J.; Lee, L.; and Danescu-Niculescu-Mizil, C. 2013. Characterizing and curating conversation threads: Expansion, focus, volume, re-entry. In *International Conference on Web Search and Data Mining (WSDM)*. ACM.
- Bakshy, E.; Rosenn, I.; Marlow, C.; and Adamic, L. 2012. The role of social networks in information diffusion. In *International Conference on World Wide Web (WWW)*, 519–528. ACM.
- Bessi, A.; Coletto, M.; Davidescu, G.; Scala, A.; Caldarelli, G.; and Quattrociocchi, W. 2015a. Science vs conspiracy:

- Collective narratives in the age of misinformation. *PLOS ONE* 10(2):1–17.
- Bessi, A.; Petroni, F.; Del Vicario, M.; Zollo, F.; Anagnostopoulos, A.; Scala, A.; Caldarelli, G.; and Quattrociocchi, W. 2015b. Viral misinformation: The role of homophily and polarization. In *International Conference on World Wide Web (WWW) Companion*. ACM.
- Borge-Holthoefer, J.; Meloni, S.; Gonçalves, B.; and Moreno, Y. 2013. Emergence of influential spreaders in modified rumor models. *Journal of Statistical Physics* 151(1-2):383–393.
- Borge-Holthoefer, J.; Rivero, A.; and Moreno, Y. 2012. Locating privileged spreaders on an online social network. *Physical Review E* 85(6):066123.
- Centola, D., and Macy, M. 2007. Complex contagions and the weakness of long ties. *American Journal of Sociology* 113(3):702–734.
- Centola, D. 2010. The spread of behavior in an online social network experiment. *Science* 329(5996):1194–1197.
- Cheng, J.; Adamic, L.; Dow, P.; Kleinberg, J.; and Leskovec, J. 2014. Can cascades be predicted? In *International Conference on World Wide Web (WWW)*. ACM.
- Cheng, J.; Adamic, L.; Kleinberg, J.; and Leskovec, J. 2016. Do cascades recur? In *International Conference on World Wide Web (WWW)*, 671–681. ACM.
- Cohen, J. 1968. Weighted kappa: Nominal scale agreement provision for scaled disagreement or partial credit. *Psychological Bulletin* 70(4):213.
- Cozzo, E.; Banos, R.; Meloni, S.; and Moreno, Y. 2013. Contact-based social contagion in multiplex networks. *Physical Review E* 88(5):050801.
- Del Vicario, M.; Bessi, A.; Zollo, F.; Petroni, F.; Scala, A.; Caldarelli, G.; Stanley, H.; and Quattrociocchi, W. 2016. The spreading of misinformation online. *Proceedings of the National Academy of Sciences (PNAS)* 113(3):554–559.
- Diao, Q.; Jiang, J.; Zhu, F.; and Lim, E. 2012. Finding bursty topics from microblogs. In *Annual Meeting of the Association for Computational Linguistics: Long Papers - Volume 1 (ACL)*, 536–544. Association for Computational Linguistics.
- Doerr, B.; Fouz, M.; and Friedrich, T. 2012. Why rumors spread so quickly in social networks. *Communications of the ACM* 55(6):70–75.
- Facebook. 2016. Company info. <http://newsroom.fb.com/company-info/>. [Online; accessed 10-January-2017].
- First Draft News Coalition. 2016. About First Draft. <https://firstdraftnews.com/about/>. [Online; accessed 10-January-2017].
- Fowler, J., and Christakis, N. 2010. Cooperative behavior cascades in human social networks. *Proceedings of the National Academy of Sciences (PNAS)* 107(12):5334–5338.
- Friggeri, A.; Adamic, L.; Eckles, D.; and Cheng, J. 2014. Rumor cascades. In *AAAI Conference on Weblogs and Social Media (ICWSM)*.
- Fronczak, A.; Fronczak, P.; and Hołyst, J. 2004. Average path length in random networks. *Physical Review E* 70(5):056110.
- Hanley, J., and McNeil, B. 1982. The meaning and use of the area under a receiver operating characteristic (roc) curve. *Radiology* 143(1):29–36.
- Haykin, S. 1998. *Neural Networks: A Comprehensive Foundation*. Prentice Hall PTR, 2nd edition.
- Ho, T. 1995. Random decision forests. In *International Conference on Document Analysis and Recognition (ICDAR)*, volume 1, 278–282. IEEE.
- Holland, P., and Leinhardt, S. 1971. Transitivity in Structural Models of Small Groups. *Small Group Research* 2(2):107–124.
- Hong, L.; Dan, O.; and Davison, B. 2011. Predicting popular messages in twitter. In *International Conference Companion on World Wide Web (WWW)*. ACM.
- Izenman, A. 2013. Linear discriminant analysis. In *Modern Multivariate Statistical Techniques*. Springer. 237–280.
- Kramer, A.; Guillory, J.; and Hancock, J. 2014. Experimental evidence of massive-scale emotional contagion through social networks. *Proceedings of the National Academy of Sciences (PNAS)* 111(24):8788–8790.
- McPherson, M.; Smith-Lovin, L.; and Cook, F. 2001. Birds of a feather: Homophily in social networks. *Annual Review of Sociology* 415–444.
- Moreno, Y.; Nekovee, M.; and Pacheco, A. 2004. Dynamics of rumor spreading in complex networks. *Physical Review E* 69(6):066130.
- Newman, N. and Levy, D.A.L. and Nielsen, R.K. 2015. Reuters Institute digital news report. <http://www.digitalnewsreport.org/survey/2015/>. [Online; accessed 10-January-2017].
- Newman, M. 2002. Assortative mixing in networks. *Physical Review Letters* 89(20):208701.
- Quattrociocchi, W.; Caldarelli, G.; and Scala, A. 2014. Opinion dynamics on interacting networks: Media competition and social influence. *Scientific Reports* 4.
- Quattrociocchi, W.; Scala, A.; and Sunstein, C. 2016. Echo chambers on facebook. Available at SSRN: <https://ssrn.com/abstract=2795110>.
- Rogers, K. and Bromwich, J.E. 2016. The hoaxes, fake news and misinformation we saw on election day. <http://www.nytimes.com/2016/11/09/us/politics/debunk-fake-news-election-day.html>. [Online; accessed 10-January-2017].
- Romero, D.; Meeder, B.; and Kleinberg, J. 2011. Differences in the mechanics of information diffusion across topics: Idioms, political hashtags, and complex contagion on twitter. In *International Conference on World Wide Web (WWW)*, 695–704. ACM.
- Salganik, M.; Dodds, P.; and Watts, D. 2006. Experimental study of inequality and unpredictability in an artificial cultural market. *Science* 311(5762):854–856.

- Seo, E.; Mohapatra, P.; and Abdelzaher, T. 2012. Identifying rumors and their sources in social networks. In *SPIE Defense, Security, and Sensing*, 83891I–83891I. International Society for Optics and Photonics.
- Silverman, C.; Strapagiel, L.; Shaban, H.; Hall, E.; and Singer-Vine, J. 2016. Hyperpartisan Facebook pages are publishing false and misleading information at an alarming rate. <https://www.buzzfeed.com/craigsilverman/partisan-fb-pages-analysis>. [Online; accessed 10-January-2017].
- Sunstein, C. 2002. The law of group polarization. *Journal of Political Philosophy* 10(2):175–195.
- Ugander, J.; Backstrom, L.; Marlow, C.; and Kleinberg, J. 2012. Structural diversity in social contagion. *Proceedings of the National Academy of Sciences (PNAS)* 109(16):5962–5966.
- Wiens, J.; Horvitz, E.; and Gutttag, J. 2012. Patient risk stratification for hospital-associated c. diff as a time-series classification task. *Advances in Neural Information Processing Systems (NIPS)* 467–475.
- Zollo, F.; Bessi, A.; Del Vicario, M.; Scala, A.; Caldarelli, G.; Shekhtman, L.; Havlin, S.; and Quattrociocchi, W. 2015a. Debunking in a world of tribes. *arXiv preprint*, <http://arxiv.org/abs/1510.04267>.
- Zollo, F.; Novak, P. K.; Del Vicario, M.; Bessi, A.; Mozetič, I.; Scala, A.; Caldarelli, G.; and Quattrociocchi, W. 2015b. Emotional dynamics in the age of misinformation. *PLOS ONE* 10:1–22.