

# HY484 Project Report

16/1/21

## CLUSTERING AND MISINFORMATION

### Motivation

The spread of misinformation has been a timeless problem. In current times though it seems that the spread speed and length have been augmented by the easiness provided by online social structures. More recently, we have seen not only incorrect information surface society, but also dangerous. Examples range from simply untenable, to hilarious claims. Such rumours have been that covid-19 virus is a made up story, and its symptoms are a direct cause of 5G antennas placed worldwide, or that the vaccines targeting it, contain microchips purposed to control the public. But hilarious as they seem, they can also have actual dangerous side effects, to the general public, such as large masses of people not cohering to social distancing, or denying to get vaccinated.

This is surely facilitated by those individual's lack of specific(or any) scientific background, required to distinguish solid from random conclusions, but that kind of background is hard to obtain for most. On the other hand many attempts are being made by social media services, to detect and label, or even remove such information, but these kind of techniques have also limitations.

Finally, from personal conversations with such individuals, we could attest that a significant percentage, support anything related to conspiracy theories or anti-establishment attitudes, by default, possibly as a way of anger outlet against unpleasant circumstances they have no control over. But, on the other hand a respectful percentage, once the truth is explained to them, feel guilty of believing and contributing to the further spread of such information, and express a yearning to avoid getting into a same predicament in the future.

### Hypothesis

As important as the aforementioned precaution measures might be, they fall on the passive side, meaning that they only address the online services side. It definitely seems beneficial to address the spreader side too. Some critical questions arise, such as what controllable user preconditions lead to them falling a victim of such behaviours, as well as how much would the misinformation cascade be mitigated by 'turning' at least some of the previous spreaders.

Firstly, we hypothesise that by 'healing' as many possible spreaders, not only the misinformation cascade is disempowered, but also, the possible valid information or response cascades, are supplemented. Secondly we hypothesise, and aim to prove, that the so called 'Emperor's Dilemma' that stands for physical social networks, also stands for the online ones. Meaning, that a cascade of an extreme claim or notion rejected by popular majority and scientific communities, can only be facilitated by individuals with ties to a small circle of people, rather than individuals belonging to a broader community. Of course a lot small

circles(clusters) like these, could finally form connections between them as they cascade, forming a bigger network of many small interconnected clusters. On the other hand the hypothesis is that individuals with a broader circle of connections, belonging to something resembling a larger fully connected component, by being exposed to more and differing opinions, develop an ability for better critical thinking, thus have less chances of falling into the unwanted behaviour.

## Methodology - Findings

For the study, two subsets of Twitter users were chosen as online social networks. These subsets were defined by users tweeting real and fake news, originating from the web, gathered and fact checked by two different sites, BuzzFeed and POLITIFACT. [1]

The sizes of the networks were 15257 users for BuzzFeed and 23865 users for POLITIFACT.

- Our first task was to determine if user behaviour was falling within discrete patterns. This was actually confirmed, what we found was that in both datasets users had a dominant tendency to post either true or false news.

For BuzzFeed the users posting only real news were 7316(48%), the users posting only fake news were 7406(48.5%), and the users posting both real/fake news were 535(3.5%).

For POLITIFACT the users posting only real news were 4437(18.6%), the users posting fake news were 19428(79%), and the users posting both real/fake news were 566(2.4%).

As we can see the cross section of users posting both types of news is very small on both networks and is not going to draw our attention from now on. Besides the users being largely split to two separate categories, we are also noticing a significantly larger percentage of fake news spreaders on POLITIFACT network.

Also it is clear at this point that both networks are homophilic in nature.

For BuzzFeed:

$$A = p^2 + q^2 = 7316^2 + 7406^2 = 108372692$$

$$B = 2 \cdot p \cdot q = 2 \cdot 7316 \cdot 7406 = 108364592$$

For POLITIFACT:

$$A = p^2 + q^2 = 4437^2 + 19428^2 = 397154153$$

$$B = 2 \cdot p \cdot q = 2 \cdot 4437 \cdot 19428 = 172404072$$

So in both cases  $A > B$ .

- Our second task was to calculate the average tweet/re-tweet rate of each group. What we found was users posting real news had marginally greater tweet rates than ones posting fake news, on both networks. Specifically, on BuzzFeed, real news posters had an average of 1.41 tweets while fake news posters had an average of 1.37 tweets.  
On POLITIFACT, real news posters had an average of 1.51 tweets while fake news posters had an average of 1.27 tweets. This could be interpreted as real news posters having a slightly more active online behaviour, while on the other hand fake news spreaders being less expressive.

---

<sup>1</sup> <https://github.com/KaiDMML/FakeNewsNet/tree/old-version>

Average re-tweet rates had insignificant differences between groups, being all very close to 1.1 retweets per tweet. This shows no persistence or urgency of opinion differentiation for any of the groups.

- With our third task we started diving into the social characteristics of the networks, calculating each group's average of user's they follow, and users they are followed by. These findings were not supportive of the hypotheses.

Real spreaders of BuzzFeed had an average of 40.6 followers and followed 40.9 other users, while fake spreaders had an average of 42.6 followers and followed 41.9 others.

Real spreaders of POLITIFACT had an average of 24.8 followers and followed 25.3 other users, while fake spreaders had an average of 24 followers and followed 23.6 others.

These averages, as well as the max followers/following recorded for each group, were found to be comparable.

- Our fourth task was to measure following activity among users of differing groups. This can be seen as a measure of diversity among groups. Metrics were taken for all four combinations of both networks. That is edges between real-to-real, real-to-fake, fake-to-real and fake-to-fake posters.

BuzzFeed results were not indicative of any differentiation among groups, since all four metrics were within 19-21 connections.

POLITIFACT users painted a slightly different picture. The real-to-real and fake-to-real connections averaged at 4.5-4.9 connections. The real-to-fake and fake-to-fake connections though, were almost four times higher, averaging at 18.4. This can be attributed on the fact that the specific network 4 times higher in fake news posters, than real.

- Moving to our fifth task, also meant moving from custom Python3 scripts, to the Gephi graph analysis/visualisation tool. This would help us further our understanding of the group's following activity. First we split the original edge lists to four new ones, each one representing one of the four categories, real→real, real→fake, fake→real, fake→fake. After the analysis, we noticed by the metrics in *table1* / *table2*, that both real and fake tweeters, network favourably with each other. This conclusion is attributed to the higher graph density plus less amount of strongly connected components and higher clustering coefficient. In social terms we are seeing the expected homophilic behaviour being expressed in both groups.
- Our sixth task was to analyse the real and fake news spreaders following activity directed to users of both groups combined. In these metrics, *table3* / *table4* we are noticing a higher clustering coefficient, in the fake news spreaders groups, but also a lower amount of strongly connected components. This in fact could be interpreted as the fake news spreaders being overall slightly more connected than the real news spreaders.

<i>TABLE<sub>1</sub></i> BuzzFeed	Degree / avg	Diameter	Path length / avg	Density	Weakly connected components	Strongly connected components	Clustering coefficient / avg
Real→real	20.38	12	3.68	0.0030	11	2121	0.09
Real→fake	12.65	1	1	0.0010	42	11750	0
Fake→real	12.55	1	1	0.0010	51	11595	0
Fake→fake	22.47	11	3.72	0.0030	15	2104	0.1320

<i>TABLE<sub>2</sub></i> PolitiFact	Degree / avg	Diameter	Path length / avg	Density	Weakly connected components	Strongly connected components	Clustering coefficient / avg
Real -> real	6.43	13	4.8	0.0020	35	1481	0.0650
Real -> fake	5.31	1	1	0	161	16534	0
Fake -> real	6.52	1	1	0	160	13092	0
Fake -> fake	19.417	14	4.19	0.0010	90	5721	0.0750

<i>TABLE<sub>3</sub></i> BuzzFeed	Degree / avg	Diameter	Path length / avg	Density	Weakly connected components	Strongly connected components	Clustering coefficient / avg
Real -> All	20.99	12	3.86	0.0020	28	9010	0.0820
Fake -> All	22.127	12	3.71	0.0020	25	8783	0.1220

<i>TABLE<sub>4</sub></i> PolitiFact	Degree / avg	Diameter	Path length / avg	Density	Weakly connected components	Strongly connected components	Clustering coefficient / avg
<b>Real -&gt; All</b>	6.302	14	4.8	0	105	15497	0.0620
<b>Fake -&gt; All</b>	19.557	14	4.18	0.0010	119	9957	0.07

## Conclusion

By studying two online social networks, with both genuine and in-genuine information, we investigated whether individuals spreading real news, have a broader circle of acquaintances, compared to the ones spreading fake news.

In case the argument held, it would be easy to follow up with some guidelines. One such example would be, that following more and/or diverse individuals on online social networks, could reduce the possibility of believing and spreading false information.

Judging by the results generated by the datasets used here, we could not confirm this claim. Most metrics along both networks used, were found either marginally different or conflicting, rendering our work inconclusive, at least up to this point.

## Future Work

It would be considered reasonably beneficial to perform similar analyses on more analogous datasets. Also the groups of users posting both real and fake news, that was not analysed in this study, due to its small relative size on both networks, could be of further interest. Although small, compared to the other two categories, it could play a pivotal role. Firstly, because of its users being possibly closer, to changing their behaviour to thorough fact checking. Secondly due to the group's possible higher social influence, if found to be more integrated to the community than the other two. Lastly a qualitative comparison, could be performed on the news subjects followed by the real versus fake spreaders. This could help us find out if having a broader circle of interests, has an effect on the probability of believing fake news or not.

## References

<https://ndg.asc.upenn.edu/wp-content/uploads/2016/04/Centola-et-al-2005-AJS.pdf>

<https://grantome.com/grant/NSF/SES-0241657>

<https://arxiv.org/pdf/2008.00791.pdf>