

A Sociolinguistic Study of Online Echo Chambers on Twitter

Nikita Duseja*

Computer Science and Engineering
Texas A&M University
nduseja@tamu.edu

Harsh Jhamtani*

School of Computer Science
Carnegie Mellon University
jharsh@cmu.edu

Abstract

Online social media platforms such as Facebook and Twitter are increasingly facing criticism for polarization of users. One particular aspect which has caught the attention of various critics is presence of users in echo chambers - a situation wherein users are exposed mostly to the opinions which are in sync with their own views. In this paper, we perform a socio-linguistic study by comparing the tweets of users in echo chambers with the tweets of users not in echo chambers with similar levels of polarity on a broad topic. Specifically, we carry out a comparative analysis of tweet structure, lexical choices, and focus issues, and provide possible explanations for the results.

1 Introduction

An echo chamber refers to a social phenomenon in which most of the content one receives in one's social media feed is heavily skewed toward one's own opinion, often defined in context of controversial political topics (Garimella et al., 2018). In social media environments, users are exposed to several polarized views on political topics. According to the selective exposure theory (Frey, 1986), individuals have a tendency to consume information from like minded individuals content and avoid contrasting perspectives. This leads to the existence of polarized segregated communities in social media, with resounding similar views. This can be concerning as such users are not exposed to alternate or opposing perspectives, which may adversely impact deliberative democratic processes.

In this work we carry out a comparative analysis of tweets from users in echo chambers versus tweets from users not in echo chambers. Specifically, we compare some properties pertaining to tweet structure, lexical choices, and top-

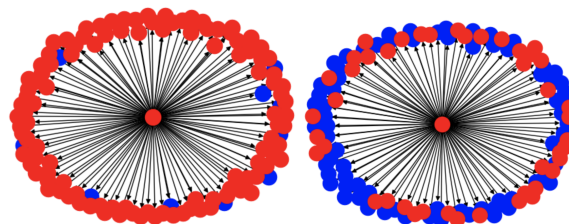


Figure 1: Visualizing user-follower network for a user in an echo chamber (center of left sub-figure) and another user not in an echo chamber (center of right sub-figure) on Twitter. A red circle represents a user with positive polarity scores for the topic *Obamacare*, while a blue circle represents a user with negative polarity scores. An arrow from x to y means that user x follows user y . A user in an echo chamber is exposed only to views very similar to his/her own opinions while a user not in an echo chamber witnesses opposing views as well. We focus on analyzing differences in tweets from the two types of users.

ics/attributes discussed in tweets. Table 1 summarizes some hypotheses of interest we framed to compare the tweets. To perform the analysis, we identify two sets of users with similar polarity levels on a topic, with one set of users being in echo chambers, while the other control set of users are not in echo chambers. We build on prior works (Garimella et al., 2018) to identify such sets of twitter users on topics such as Affordable Care Act (*Obamacare*), a comprehensive health care reform law that was enacted under the Obama administration. Figure 1 provides a pictorial overview of the network of a user in an echo chamber and another user not in an echo chamber.

There have been many recent works focusing on echo chambers in online social media (Garimella et al., 2018; Grömping, 2014; Barberá et al., 2015; Kwon et al., 2012). Many prior works point out presence of a large number of social media users in echo chambers (An et al., 2014; Bakshy et al., 2015; Lawrence et al., 2010). Our primary contribution is a comparative analysis of tweets from users in echo chambers versus users not in echo

* ND and HJ contributed equally for this paper

chambers. We build on work of [Garimella et al. \(2018\)](#), who focus on identifying users in echo chambers and characterizing various social network properties of such users. We believe our study can help in further characterizing and understanding online echo chambers, which may help in mitigating negative impacts associated with echo chambers.

2 Dataset

We use data from [Garimella et al. \(2018\)](#) who calculated polarity of twitter users towards topics such as *Obamacare* on twitter. The dataset contains polarities of users, as well as the user-follower graph for the same set of users.

We choose to work with *Obamacare* topic for the analysis in this paper. All the polarity scores of users are in the range $(-2.5, 2.5)$. A higher positive score represents more *conservative* viewpoint, while a more negative score represents a more *liberal* viewpoint ([Garimella et al., 2018](#)) (Not to be confused with sentiment towards the topic). In general, positive score users can be considered as *conservative* users, while negative score users can be considered as *liberal* users.

The dataset also consists of user-follower network in Twitter for the relevant set of users. We are mainly interested in finding the followees of users since a user is typically exposed to tweets and re-tweets of his/her followees in the social media feed.

2.1 Echo chambers

We consider following notation for a twitter user $u \in U$, where U is the set of all users under consideration.

- $S(u)$: Set of followees of u having same polarity as the user u
- $D(u)$ Set of followees of u having different polarity as the user u

To characterize an echo chamber, we define *Homophily score* $H(u)$ for a user u as follows:

$$H(u) = \frac{|S(u)| - |D(u)|}{|S(u)| + |D(u)|} \quad (1)$$

Thus, $H(u) \in (-1, 1)$ range (both inclusive). A score of $H(u) = 1$ means that all the followees of the user u have same polarity about the given topic as u himself/herself. This characterizes an

extreme case of being in an echo chamber. On the other hand, score close to 0 would suggest that the user has followees of both polarity, i.e they belong to classes on both sides of the spectrum in equal number. We define a threshold of θ_1 such that users with $H(u)$ above θ_1 are said to be in echo chamber **EC**. Users having $H(u) < \theta_2$ ($\theta_2 < \theta_1$) are said to not be in an echo chamber **NE**. We first report results for $\theta_1 = 0.9$, $\theta_2 = 0.7$, and later discuss the robustness to these choices.

$$EC = \{u : u \in U, H(u) \geq \theta_1\} \quad (2)$$

$$NE = \{u : u \in U, H(u) < \theta_2\} \quad (3)$$

3 Methodology

Our aim is to compare tweets from users in echo chambers against tweets of users not in echo chambers. Towards this end, we control for the polarity of the users, so that we could study any differences in nature of tweets which correlates with being in an echo chamber. Specifically, we restrict users to polarity range 0.5 to 1.5. We term users in this polarity region as *moderate conservatives* (MC).

We work with moderate conservatives and compare the tweets of users in echo chambers versus users not in echo chambers. Next, we filter to retain only those tweets which talk about *Obamacare* (the original data contains all tweets of users over a long time duration). This leaves us with 47,533 tweets for moderate conservatives in echo chambers **MC-EC**, and 35,820 tweets of moderate conservatives not in echo chamber **MC-NE**, talking about *Obamacare*.

Table 1 lists our main hypotheses which we test on the dataset. We use chi-squared tests for testing statistical significance while comparing counts of features. The chi-squared test ([Greenwood and Nikulin, 1996](#)) is used to determine whether a perceived association between two categorical variables is by chance or reflects a real association between these two variables in the data. It compares the observed frequency with expected frequency (expected assuming no correlation).

We perform some pre-processing on the tweet texts before conducting the experiments. We use Twitter tokenizer from NLTK ([Bird et al., 2009](#)) to tokenize the tweets. We retain the hashtags and URLs as they are needed in the experiments and analysis.

Type	Hypothesis	Holds in data?
Tweet-structure	MC-NE tweets are more likely than MC-EC to cite external resource	Yes
Tweet-structure	MC-EC tweets are more likely than MC-NE to contain hashtags	Yes
Vocabulary	MC-NE tweets are more likely than MC-EC to express uncertainty	Yes
Vocabulary	MC-NE tweets are more likely than MC-EC to use swear words	Yes
Topical	Certain topics are talked about more in MC-EC tweets and vice versa	Yes

Table 1: Summary of the hypotheses and the results. We carry out tweet structure analysis, vocabulary choice analysis, and a topic-level analysis, and observe significant difference in the tweets from the two types of users.

4 Experiments

In this section, we describe more details about the experiments and corresponding observations. We define three types of hypotheses: 1) Tweet Structure Analysis 2) Vocabulary Analysis 3) Topic Analysis.

4.1 Tweet structure analysis

Tweet structure analysis aims to uncover differences in tweets from the two types of users with regard to aspects like use of hash-tags, use of accompanying URLs, etc. Our first hypothesis is based on intuition that MC-EC users feel less compelled to cite an external resource while tweeting, as all the tweets in their feed already resonate with their own view-points. Our second hypothesis is that MC-EC tweets are more likely to contain hash-tags, as MC-EC users may more strongly believe in correctness of their own viewpoint, and may use more hash-tags with the intention to spread their strongly believed view-points.

4.1.1 Evidence / Link citing

Hypothesis: We hypothesize that users not in echo chamber may be more likely to tweet or re-tweet with citing external news link or other sources. This follows the general idea that users in an ‘echo chamber’ might feel less of a need to justify their claims or opinions as people (followers) around them *echo* with similar opinions.

Analysis: We perform the test using a chi-squared test. We first identify URLs using simple regular expressions. We notice that most of the urls were shortened URLs. We leverage python library BeautifulSoup (Richardson, 2007) to identify the expanded URLs, and then filter out any twitter URLs (since these often correspond to other user’s status’). We observe that tweets from MC-NE users are much more likely to contain external URLs, which are often news or opinion pieces,

compared to users in echo chambers($p < 0.01$ as per chi-squared test). Specifically, about 35% tweets from MC-NE users had an external link while only about 19% of the tweets from MC-EC contained an external link.

4.1.2 Use of hashtag

Hypothesis: Hashtags are widely used in tweets, often to explicitly tag the tweet about being a specific topic or point, and are often used with the intention to spread messages or viewpoints. We hypothesize that MC-EC users’ tweets are more likely to use hashtags.

Analysis: We explore the degree to which hash-tags are used in tweets between the two types of users. We observe that our hypothesis holds in the dataset - tweets and re-tweets from MC-EC users are much more likely to contain hashtags($p < 0.01$ as per chi-squared test). Specifically, 45909 MC-EC tweets out of 47533 had at least one hashtag, while only 17097 MC-NE tweets out of 35820 had hashtags.

4.2 Vocabulary analysis

We inspect if being in an echo chamber is correlated with more/less usage of specific types of words. For example, MC-NE users are exposed to varying viewpoints, and therefore their vocabulary choice might reflect some uncertainty in views.

4.2.1 Words expressing uncertainty

Hypothesis: Since users in echo chambers are exposed only to opinions similar to theirs in the online media, they might show more certainty in their tweets. Similarly, users not in echo chambers are exposed to alternative views also in online social media, and as such may use uncertainty expressing words more frequently. We hypothesize that tweets from MC-NE users are more likely to contain uncertainty depicting words.

We use following list of uncertainty depicting words: ‘may’, ‘might’, ‘perhaps’, ‘maybe/may-be’, ‘possibly’, ‘likely’.

Analysis: We test the hypothesis using chi-squared test and observe that the usage is more frequent in users outside of echo chambers ($p < 0.01$). For example, word ‘might’ appears 238 times in MC-NE tweets, while occurs only 159 times in MC-EC tweets. Word ‘may’ appears 720 times in MC-NE, while occurs only 581 times in MC-EC.

4.2.2 Use of swear words

Hypothesis: We hypothesize that users not in echo chambers are more likely to express frustration through swear words on witnessing opposing viewpoints. We use the list of common English swear words¹. We do expand the list to include commonly occurring variants. For example, *f**cking* is a commonly used word which would not have shown up on doing exact token match to *f**ck*. This expansion was done manually, as automatic lemmatization tools did not work satisfactorily. Improving swear word detection would be part of future work.

Analysis: We inspect the total count of swear words used in the two set of tweets. We observe that the proposed hypothesis holds in the data ($p < 0.01$ using a chi-squared test). For example, word ‘f**k’ appears 41 times in MC-NE tweets, while occurs only 1 time in our set of MC-EC tweets. This analysis suggests that being in echo chamber is correlated with lesser use of swear words.

4.3 Topic analysis

We had filtered the tweets to be about the broad topic of Obamacare. In this analysis, we are interested in comparing the main (sub-)topics about Obamacare that are discussed in the two user groups. Towards this goal, we run a topic model on the set of tweets from each user group.

Hypothesis: We hypothesize that certain topics would be correlated with presence in echo chamber i.e. some topics would be talked about more in MC-EC tweets while certain other topics would be covered more in MC-NE tweets. We believe that presence in echo chambers might have an

¹ Available at https://en.wiktionary.org/wiki/Category:English_swear_words

effect on the aspects of *Obamacare* that users are tweeting about.

Analysis: We run LDA (Latent Dirichlet Allocation) topic model to extract the topics for the combined set of tweets from both user types. LDA is a generative model, where each ‘document’ is supposed to have been generated using a multinomial distribution over the set of topics, and each word in the document can be thought of being generated from a topic picked up based on the drawn topic distribution for that document. Each topic itself is a multinomial distribution over the vocabulary of words.

Towards this end, we run the topic model for K number of topics on all the tweets i.e. tweets from MC-NE and MC-EC combined. We limit the vocabulary to 1000 most frequently occurring words in the dataset excluding the English stop words. The model approximates the multinomial distribution over vocabulary for each topic, and also computes the relative proportion of each topic for every *document* (tweet). For each of the two types of tweets (from MC-NE and MC-EC users), we compute an aggregated topic distribution for that type by summing the topic distribution vectors of corresponding tweets. This can be thought of summing up fractional counts of occurrence of topics, and this provides us with two topic occurrence pseudo-counts, one for each set of tweets.

For each topic, we test if it’s pseudo count is significantly different between the two types of tweets. We experiment with $K = 10, 15, 20$, and observe that 11 (for $K = 20$), 6 (for $K = 10$) and 3 (for $K = 5$) topics had statistically significant different occurrence in the two sets of tweets ($p < 0.01$). For example, for $K = 20$, topic *future* (‘time’, ‘year’, ‘watch’, ‘future’) occurs much more in MC-EC tweets while topic *repeal* (‘repeal’, ‘vote’, ‘repeal’, ‘senate’) is present much more in MC-NE.

4.4 Sensitivity to homophily threshold parameters θ_1 and θ_2

Above experiments were reported for homophily thresholds $\theta_1 = 0.9, \theta_2 = 0.7$. These values were selected such that users with very high homophily scores are marked as being in echo chambers. We repeat the all the experiments with more sets of parameter values ($\theta_1 = 0.8, \theta_2 = 0.8$) and ($\theta_1 = 0.9, \theta_2 = 0.9$), and observe very sim-

ilar results, with all the tested hypotheses leading to same conclusions. This demonstrates that the analysis is robust to changes in these parameter values.

5 Discussions

Since we were interested in comparing the linguistic properties of the tweets from two types of users, we control for polarity levels, and select MC-EC and MC-NE sets of users. It is possible to extend the work on other such pairs of categories, such as ML-EC and ML-NE (ML: moderate liberals), and test the generality of the proposed hypotheses to other such groups. Moreover, we experiment with tweets on only one topic: Obamacare - it would be interesting to test the generalizability of the hypotheses to other data sets as well. Such extensions to the current work are part of future directions.

There are certain limitations of the current analysis. We did not take into account many network and content popularity effects. Moreover, we do not comment on any causality aspect: for example, does one's presence in an echo chamber makes one's tweet less likely to contain uncertainty depicting words, or if less polarized users are less likely to get trapped in an echo chamber. This remains an important possible future extension.

Related Work Garimella et al. (2018) propose methods to identify partisan and bipartisan users, and characterize such users based on network effect, profile information, and interaction actions such as retweets. We leverage their work and dataset to define echo chambers. However, our main focus is to do a linguistic comparison of tweets based on whether the tweet is from a user in an echo chamber or not.

Many prior works (Garimella et al., 2018) aim to study echo chambers in context of various network effects. Some prior works correlate retweet network with political ideology of the users (Barberá et al., 2015). Bakshy et al. (2015) study the consumption of online content generated by users of opposing views. Gilbert et al. (2009) conduct a comment based study on political blogs and find that, to a great extent, comments are in agreement with the views of the author of the blog. We on the other hand correlate some linguistic properties of tweets with presence in an online echo chamber.

6 Conclusion

We have carried out a comparison of tweets between moderate conservatives in echo chambers with moderate conservatives not in echo chambers, on the topic *Obamacare*. We carry out analysis at three different levels: tweet-structure level, topic level, and word-group level. We observe statistically significant difference in frequency of usage of uncertainty depicting words, hashtags, swear words, and external URL links, as well as a difference in the aspects of *Obamacare* talked about frequently between the two types of tweets. We also highlight possible future extensions to our work.

Acknowledgements

We are thankful to David R. Mortensen for insightful comments and discussions. We also acknowledge Lori Levin, Vidhisha Balachandran, Sanket V. Mehta, Kundan Krishna, and anonymous workshop reviewers for providing valuable feedback. We are thankful to Kiran Garimella for sharing dataset with us.

References

- Jisun An, Daniele Quercia, Meeyoung Cha, Krishna Gummadi, and Jon Crowcroft. 2014. Sharing political news: the balancing act of intimacy and socialization in selective exposure. *EPJ Data Science*, 3(1):12.
- Eytan Bakshy, Solomon Messing, and Lada A Adamic. 2015. Exposure to ideologically diverse news and opinion on facebook. *Science*, 348(6239):1130–1132.
- Pablo Barberá, John T Jost, Jonathan Nagler, Joshua A Tucker, and Richard Bonneau. 2015. Tweeting from left to right: Is online political communication more than an echo chamber? *Psychological science*, 26(10):1531–1542.
- Steven Bird, Ewan Klein, and Edward Loper. 2009. *Natural language processing with Python: analyzing text with the natural language toolkit*. "O'Reilly Media, Inc."
- Dieter Frey. 1986. Recent research on selective exposure to information. In *Advances in experimental social psychology*, volume 19, pages 41–80. Elsevier.
- Kiran Garimella, Gianmarco De Francisci Morales, Aristides Gionis, and Michael Mathioudakis. 2018. Political discourse on social media: Echo chambers, gatekeepers, and the price of bipartisanship. In *Proceedings of the 2018 World Wide Web Conference*

- on *World Wide Web*, pages 913–922. International World Wide Web Conferences Steering Committee.
- Eric Gilbert, Tony Bergstrom, and Karrie Karahalios. 2009. Blogs are echo chambers: Blogs are echo chambers. In *2009 42nd Hawaii International Conference on System Sciences*, pages 1–10. IEEE.
- Priscilla E Greenwood and Michael S Nikulin. 1996. *A guide to chi-squared testing*, volume 280. John Wiley & Sons.
- Max Grömping. 2014. echo chambers partisan facebook groups during the 2014 thai election. *Asia Pacific Media Educator*, 24(1):39–59.
- K Hazel Kwon, Onook Oh, Manish Agrawal, and H Raghav Rao. 2012. Audience gatekeeping in the twitter service: An investigation of tweets about the 2009 gaza conflict. *AIS Transactions on Human-Computer Interaction*, 4(4):212–229.
- Eric Lawrence, John Sides, and Henry Farrell. 2010. Self-segregation or deliberation? blog readership, participation, and polarization in american politics. *Perspectives on Politics*, 8(1):141–157.
- Leonard Richardson. 2007. Beautiful soup documentation. *April*.