# Evidence of Demographic rather than Ideological Segregation in News Discussion on Reddit

*Keywords: Homophily, Socio-demographic factors, Polarization, News, Reddit*

## Introduction

Socializing across ideological boundaries—a phenomenon known as *affective polarization*—has been reported to be increasingly difficult, especially in the United States [1]. The causes of polarization are widely debated. Some studies have pointed at social media [2]: recommendation algorithms, algorithmic filter bubbles, and self-sorting could reinforce ideological separation. An alternative explanation lies in demographic factors [3]: underlying material rifts would make gender, age, and wealth boundaries increasingly divisive. To assess these competing hypotheses, we investigate whether online interactions are segregated by demographic boundaries, or by ideological echo chambers. In particular, we focus on one of the main loci where opinions are formed and challenged: the discussion of news [4].

## Data

We analyze interactions on the `r/news` subreddit. For each year between 2016 and 2020 (included), we select non-bot users of this subreddit above a given threshold of activity. Then, we build their interaction graph: an arc $u \to v$ represents a user $u$ replying to a message posted by user $v$. We obtain 21 to 34 thousands nodes per year, and average degree 36.8 to 41.7. Then, we extract the socio-demographic features of these users through their participation to the other Reddit communities. Waller et al. [5] assign a score to each subreddit for three demographic axes—age, gender, and affluence—and for political leaning. For each axis, we compute the score of a user as the weighted average of the scores of the subreddits they participate in. Then, we binarize such scores by labeling only users in the highest and lowest quartile for each axis, obtaining the following set $F$ of features: Young, Old, Male, Female, Poor, Rich, Left-leaning, and Right-leaning. Each user $u$ is thus associated to a binary feature vector $x_u \in \{0,1\}^{|F|}$.

## Analysis

*Socio-demographic features model.* On this data, we apply our first logistic regression model, inspired by the feature-feature graph model [6]. This model can be summarized as $\text{logit}(y_{u,v}) = \beta_0 + x_u^\mathsf{T} W x_v$, where the dependent variable $y_{u,v}$ indicates whether $u$ interacted with $v$; the independent variables $x_u \otimes x_v$ are the $|F| \cdot |F|$ variables such that the element $h,k$ is 1 if node $u$ has feature $h$ and node $v$ has feature $k$; and $W \in \mathbb{R}^{|F| \times |F|}$ is a feature-feature matrix [6] where each element $W_{h,k}$ indicates the log odds ratio for nodes with feature $h$ to interact with nodes with feature $k$. Then, we fit the model on a balanced sequence of pairs [7], where the positives are all the arcs $E$, and the negatives are a sample of non-connected pairs of nodes, drawn proportionally to their degree. We obtain an estimate of the feature-feature matrix $W$ for each year.

    Figure 1a depicts the matrix of one year (2016), but showing only the coefficients that are statistically significant in at least four years out of five. The most striking characteristic is the diagonal: left-leaning and right-leaning users interact significantly more across the boundary, and significantly less within their respective group, w.r.t. the null model. Conversely, regarding age, gender, and affluence, users tend to interact significantly more *within* their demographic

group. In other words, we observe *segregation* along demographic boundaries, in the sense of "restriction of contacts between various groups" [8].

*Socio-demographic model with topics (SD+T).* We then study whether some users are more likely to interact because they are drawn towards the same topics, and whether there are any residual effects of the socio-demographic features. To do so, we associate each interaction $(u,v)$ to a topic $t \in T$ by applying a state-of-the-art topic classifier to the title of the submission under which the interaction happens. We use these topics in a second model, summarized as $\text{logit}(y_{u,v}) = \beta_0 + x_u^\top W x_v + x_u^\top Q e_t + x_v^\top Q e_t$, where $e_t$ is the one-hot vector for $t$ and $Q \in \mathbb{R}^{|F| \times |T|}$ is a feature-topic matrix representing the log odds ratio that a user with feature $h$ comments on topic $t$.

Controlling for the topic does not alter the results of the model, although it considerably improves its fit. Figure 1b shows the weights of the matrix $W$ on the block-diagonal for each socio-demographic axis, confirming the previous observations. Figure 2 instead shows the weights of the feature-topic matrix $Q$ for 2016 (the other years are qualitatively similar). We observe meaningful associations: e.g., richer users are interested in technology, travel, business, healthy living, and entertainment.

## Discussion

We have quantitatively compared two competing hypotheses to explain affective polarization and found that demographic segregation is more prominent than ideological echo chambers on Reddit news discussions. The interaction network of `r/news` displays clear status homophily, whereby users similar in age, gender, and affluence, are more likely to interact; and also find value heterophily, as users with *opposing* political leaning are also more likely to interact. Importantly, these patterns are robust to controlling for the interest in the topic of the news. Overall, our results show that ideological echo chambers are not evident in news discussions, while we find evidence of demographic segregation. Strikingly, this demographic segregation happens even though the demographic traits at play are not immediately available to the discussants. The most likely explanation for this result supports the theory that one's worldview is affected by the demographic group they belong to.

## References

[1] S. Iyengar, G. Sood, Y. Lelkes, *Public opinion quarterly* **76**, 405 (2012). (Cited on 1)

[2] R. K. Garrett, *Journal of Computer-Mediated Communication* **14**, 265 (2009). (Cited on 1)

[3] M. Bradley, S. Chauchard, *Frontiers in Political Science* p. 64 (2022). (Cited on 1)

[4] J. Kim, R. O. Wyatt, E. Katz, *Political communication* **16**, 361 (1999). (Cited on 1)

[5] I. Waller, A. Anderson, *Nature* **600**, 264 (2021). (Cited on 1)

[6] C. Monti, P. Boldi, *Internet Mathematics* **1**, 2640 (2017). (Cited on 1)

[7] G. De Francisci Morales, C. Monti, M. Starnini, *Scientific Reports* **11**, 2818 (2021). (Cited on 1)

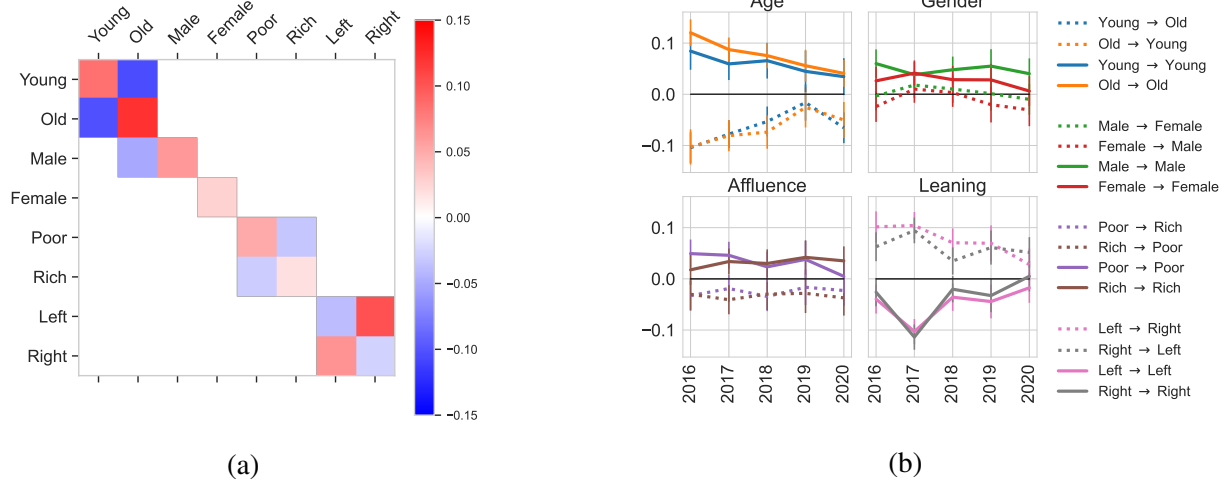[8] L. C. Freeman, *Sociological Methods & Research* **6**, 411 (1978). (Cited on 2)

Figure 1: Logistic regression coefficients for interactions between demographic features on the `r/news` subreddit. Figure (a) depicts the full feature-feature matrix obtained on 2016 data by the socio-demographic features model, but showing only coefficients that are significant in at least 4 years out of 5. Figure (b) focuses on the heterophilic and homophilic interactions, and show the associated coefficients obtained independently for each year by the SD+T model, with 95% CI.
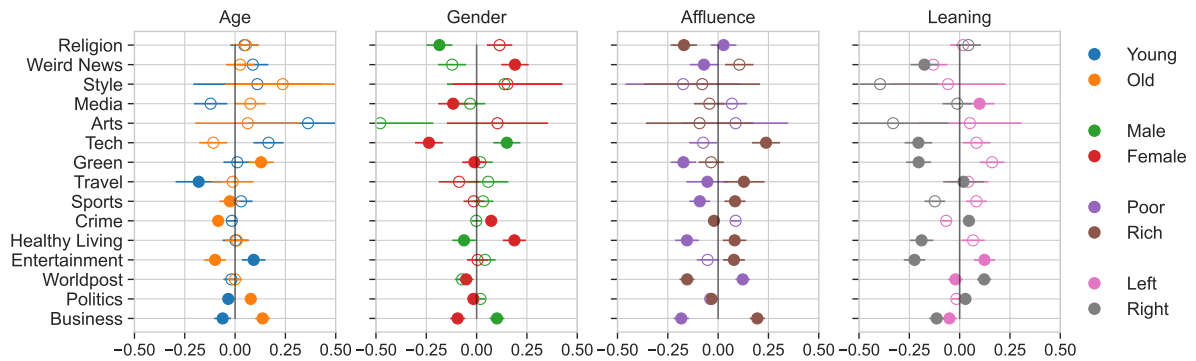


Figure 2: Logistic regression coefficients, with 95% CI, for interactions between demographic features and topics obtained for the `r/news` subreddit in 2016 by the SD+T model. Filled disks indicate associations significant in at least 4 years out of 5.