
An Author Network to Classify Open Online Discussions

Mattias Mano

i3-CRG Ecole Polytechnique,
Université Paris-Saclay
Palaiseau, France
mattias.mano@polytechnique.edu

Jean-Michel Dalle

UPMC, i3-CRG Ecole
Polytechnique
Paris, France
jean-michel.dalle@upmc.fr

Joanna Tomasik

LRI, CentraleSupélec, Université
Paris-Saclay
Orsay, France
joanna.tomasik@lri.fr

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the Owner/Author.

OpenSym '17, August 23-25, 2017, Galway, Ireland
© 2017 Copyright is held by the owner/author(s).
ACM ISBN 978-1-4503-5187-4/17/08.
<https://doi.org/10.1145/3125433.3125450>

Abstract

Among other modalities, online coordination can notably rely on discussions and forums. However, and notwithstanding increasing research efforts, direct approaches that would help communities and moderators distinguish between gossip and serious debates are still largely missing. We present an innovative methodology to detect the different structures of online discussions in the sub-Reddit *Change My View*. Applying a clustering algorithm to the author networks, we highlight three distinct classes characterized by alternative behaviors. To better understand the underlying social dynamics, we implement a relational event model that provides evidence for three effects whose influence can affect the structure of online discussions.

Author Keywords

Online discussions ; classification ; author networks ; graph motif analysis ; relational event model

ACM Classification Keywords

E.1 [Data Structures]: Graphs and Networks; I.5.3 [Pattern Recognition]: Clustering — Algorithms, Similarity measures; I.2.8 [Artificial Intelligence]: Problem Solving, Control Methods, and Search — Graph and Tree search strategies

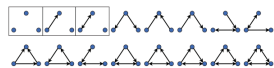


Figure 1: Triad dictionary [3]

Introduction

Open Collaboration research is a major field where online collaboration has emerged and developed on the Internet. The typical approach consists in modeling forums as graphs. In this study, we focus on the author network, modeling the link between the participants. In this case, nodes are authors and two authors are linked if one answers the other. We analyze motifs [6] *i.e.* sub-graph of three nodes in those author networks (Figure 1). We then apply the Relational Event Model — REM [1] which analyzes the dynamics behind an event stream, here, the sequential chain of posts.

METHODOLOGICAL FRAMEWORK

Reddit Change My View — CMV

We select a sub-Reddit, called *Change My View — CMV*, which admits a complex structure of discussion. It allows participants to answer any comment, with no restriction on depth of sub-comments. Such as *Agora* in Ancient Greece, someone — the Original Poster (OP) — opens a discussion announcing her idea on a topic. To do so, the OP agrees to follow the conversation at least up to three hours after the beginning and she *must personally hold the view and be willing to have it changed*¹. From there, the OP as well as any other participant, might highlight any argument that allowed her to make a step in the change of her view. Then, she attributes a delta Δ and has to explain why this comment has been convincing, even if the change is minor: a Δ does not terminate the conversation.

We work on an open database from [7], including all discussions from January 2013 (the creation of CMV) to April 2015: 15 106 discussions, composed of 1 027 074 posts, 10 470 OP and 67 693 participants. Moreover, we access to an important quantity of information, such as the number of

Δ earned by a participant, used as an experience level, or whether a post was awarded with a Δ .

Author Networks and Motifs

We extract the author network from each discussion. In these author networks, vertices represent authors. When author *A* answers author *B*, we create an arc from *A* to *B*. We count motifs in author networks, *i.e.* the number of sub-graphs composed of three nodes. As Milo *et al.* [6], we only analyzed motifs with exactly three vertices and at least two edges (Figure 1).

In order to analyze the motifs, we use a specific frequency measure: for each discussion we normalize the percentage of a given motif by its percentage in the entire database, puts formally in equation (1), with *i* for the motif, *j* for the network and $\text{motif}_{i,j}$ the number of motifs *i* in the network *j*. This ratio allows us to compare the significance of motifs in a discussion with respect to its significance in the whole dataset [2]. We compute this ratio for all 13 studied motifs, for every author network in the database and thus have a set of 13 ratios for every discussion.

DBSCAN Clustering Algorithm

The density-based clustering application with noise — DBSCAN [4] is a data clustering algorithm based on the density metrics, which does not require the number of clusters given in advance. It is configured with a distance parameter ϵ which defines the maximal distance between two points considered as neighbors and a threshold *MinPts* understood as a minimum number of points in a neighborhood.

Random Effect Relation Event Model — REREM

REM, detailed in [1], allows a more precise analysis of the discussion dynamics. It relies upon the derivation of likelihood. Having the exact timing of each post in discussions, we apply the interval model and focus on three specific ef-

$$\text{motif}_{i,j_{norm}} = \frac{\text{motif}_{i,j}}{\sum_i \text{motif}_{i,j}} \frac{\sum_j \text{motif}_{i,j}}{\sum_{i,j} \text{motif}_{i,j}} \quad (1)$$

¹<https://www.reddit.com/r/changemyview/wiki/rules>

fects [1]: *recency* — measuring a preference to answer those who recently exchange with oneself; *persistence* — measuring a social "inertia"; *participation shifts* (P-shift) — measuring a shift in authoring between two consecutive events, highlighting a tendency to local reciprocity in discussion.

Besides, due to a lack of computational power, researchers are not able to compute the REM on a large event stream, such as our database. Following [5], we apply a Random Effect REM — REREM to estimate the average effect size of REM per cluster.

RESULTS

Clustering

We compute for all networks, 13 ratios motif $_{i,j_{norm}}$ defined in the previous section. We use this set as input for the DBSCAN algorithm, selecting $MinPts = 15$ and $\varepsilon = 0.12$. DBSCAN produced three clusters which group 11 193, 2 122 and 1 427 discussions, respectively. Outliers (364, 2.4%) are negligible.

We analyze the occurrence of sub-graphs in each author network, composed of three vertices, modeling the authors, and at least two edges, modeling a discussion between two authors. Motifs $A \rightarrow B \leftarrow C$, $A \leftrightarrow B \leftarrow C$ and $A \leftrightarrow B \leftrightarrow C$ embody around 90% of motifs distribution in every cluster (Table 1). Moreover, in 98% of discussions, the opening post has the highest degree and the OP embodies in average 36.79% of the normalized betweenness centrality. Those results are an evidence in favor of OP predominant position. Thus, we hypothesize that the central node B in those particular triads is the OP the most often. With this assumption, the first and second motifs characterize a discussion where the OP answers few comments, whereas the third one characterizes threads where the OP discusses several times with her challengers.

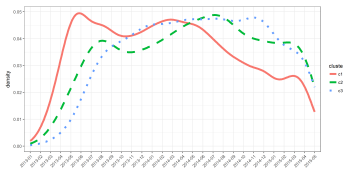


Figure 2: Probability density function of the time of the last post, per cluster: solid = cluster 1; dashed = cluster 2; dotted = cluster 3.

	C1	C2	C3
$A \rightarrow B \leftarrow C$	34.55	29.57	37.51
$A \leftrightarrow B \leftarrow C$	34.41	34.65	31.50
$A \leftrightarrow B \leftrightarrow C$	22.44	25.44	16.62

Table 1: Percent of motif by clusters

Thus, in the cluster 3, the OP has a tendency to answer less than in other clusters. Either she does not answer ($A \rightarrow B \leftarrow C$ at 37.51%, highest score across clusters) or focuses on few challengers ($A \leftrightarrow B \leftarrow C$ at 31.50%, lowest score across clusters) rather than answers all participants ($A \leftrightarrow B \leftrightarrow C$ at 16.62%, lowest score across clusters). Moreover, discussions in this cluster attract much more participants (106 in average), who comment more often (202 posts in average) and discuss longer (almost 900 hours in average).

On the contrary, cluster 2 groups discussions where the OP is more active (highest percent across clusters for $A \leftrightarrow B \leftrightarrow C$ and $A \leftrightarrow B \leftarrow C$, lowest for $A \rightarrow B \leftarrow C$ with, respectively, 25.44, 34.65 and 29.57 percent). Besides, discussions in this cluster are smaller than in the previous one with 35 authors posting 97 comments discussing over 336 hours, in average.

Finally, the cluster 1 is the most heterogeneous, grouping almost 75% of discussions, with almost as much $A \rightarrow B \leftarrow C$ as $A \leftrightarrow B \leftarrow C$ contributing for two third of the motif distribution. Moreover, it contains the smallest and shortest discussions: 15 authors write 33 messages over 144 hours in average.

Besides, we analyze the probability density function of the date of the last post (Figure 2), per cluster, as a proxy of the age of the discussion. The homogeneity of the density across clusters confirms that DBSCAN does not group by age of discussion. Thus, every difference we have highlighted hold independently of the age of the discussions: if threads in cluster 1 are shorter, it is due to the fact that authors stop discussing, not because it does not get the time to grow.

REREM

We apply a REREM on samples of clusters 2 and 3 (905 and 510 discussions respectively for cluster 2 and cluster 3), due to limited power, to put into evidence the different social mechanisms underlying them. We applied the model, with two covariates — being or not the OP and being or not an expert, defined as someone who has already received at least one Δ since her enrollment — for each of the three studied effects, one after another: recency, persistence, and P-shift effects.

First, concerning the covariates effects: being the OP or an expert implies a higher activity, even if being the OP is far more prominent than being an expert. The coefficients of the covariates stay approximately at a same level across the model but are higher in cluster 2 than in cluster 3. Thus, the OP has a higher impact in cluster 2 than in cluster 3.

Second, concerning the social dynamics, the main effect in both clusters is the local reciprocity corresponding to the P-shift effect. Still, both the recency and persistence effects are higher in the cluster 3: it accentuates the existence of plural discussions in one thread.

In a nutshell, covariates influence more the dynamics in cluster 2, with a predominant impact of being the OP, whereas in cluster 3, covariates are less significant than in cluster 2, but the social effects impact more the dynamics of discussions.

CONCLUSION

Our results suggest that, in Reddit — *Change My View*, three groups of discussions may be observed. The distinction is due firstly to the implication of the Original Poster (OP). Either, she manages to follow her challengers, building a sub-discussion with the majority of them or the discussion escapes from the OP lead and participants exchange between themselves. The third class contains discussions which do

not attract enough participants. Besides, REM highlighted different social mechanisms between the first two clusters, confirming the major impact of the OP in the first one, whereas social dynamics such as recency and persistence affect more the second one.

REFERENCES

1. C. T. Butts. 2008. 4. A Relational Event Framework for Social Action. *Sociological Methodology* 38, 1 (Aug. 2008).
2. R. Charbey and C. Prieur. 2016. Applying a structural typology on a large dataset of inline personal networks. In *Second European Conference on Social Networks*. Paris, France.
3. P. Cunningham, M. Harrigan, G. Wu, and D. O’Callaghan. 2013. Characterizing ego-networks using motifs. *Network Science* 1, 02 (Aug. 2013).
4. M. Ester, H.-P. Kriegel, J. Sander, and X. Xu. 1996. A density-based algorithm for discovering clusters in large spatial databases with noise. AAAI Press.
5. H. Liang. 2014. The Organizational Principles of Online Political Discussion: A Relational Event Stream Model for Analysis of Web Forum Deliberation. *Human Communication Research* 40, 4 (Oct. 2014).
6. R. Milo, S. Shen-Orr, S. Itzkovitz, N. Kashtan, D. Chklovskii, and U. Alon. 2002. Network Motifs: Simple Building Blocks of Complex Networks. *Science* 298, 5594 (Oct. 2002).
7. C. Tan, V. Niculae, C. Danescu-Niculescu-Mizil, and L. Lee. 2016. Winning Arguments: Interaction Dynamics and Persuasion Strategies in Good-faith Online Discussions. *arXiv:1602.01103 [physics]* (Feb. 2016).