

You are the way you structurally talk: structural-temporal neighbourhoods of posts to characterize users in online forums

(Draft version: April 17, 2016)

Alberto Lumbreras
Julien Velcin
Laboratoire ERIC
Université de Lyon, France
alberto.lumbreras@univ-lyon2.fr
julien.velcin@univ-lyon2.fr

Bertrand Jouve
FRAMESPA/IMT
Université de Toulouse 2, France
jouve@univ-tlse2.fr

Marie Guégan
Technicolor, France
marie.guegan@technicolor.com

Abstract—Users of social networks are often characterised by extracting some relevant features from the graphs associated to that network. The structural dynamic of the conversation are usually forgotten due to a lack of proper tools. We present a purely graph-based method to cluster users in online forums.

I. INTRODUCTION

The interactions between users in online forums are often modelled as complex networks. As such, they can be studied from many levels, each of which is represented by a graph where vertices and edges may represent different things. The most studied graph is a graph where vertices represent users and edges represent interactions between users (a post from one user replying to a post from other user). From the point of view of the community, some typically analysed properties are the degree distribution, the clustering coefficient, density, or diameter of the graph. From the point of view of the individual users, analyses are typically focused either on the centrality of users, on the community structure or in blockmodeling (groups of users that tend to interact with the same other groups) [1]. Another common graph is a tree graph where vertices represent posts (or e-mails) and edges from one vertex to another indicates that the first is a reply to the second. Given the different discussion trees of a forum (a forest) the global analyses include the depth distribution, branching factors, or even the time between two posts [2]. representation include depth distribution of discussions, branching factors and so. Studying discussion trees from the point of view of the users means, since trees are discussions, studying the user at a conversational level. Unfortunately, there is a lack of structural tools to perform this analysis. In this paper, we propose to fill this gap through the concept of *structural-temporal neighbourhood*.

Our goal in this paper is two-fold: on the one hand, to illustrate how structural-temporal neighbourhoods can play the role of triads for conversation trees, in the sense that they show us the local structures or dynamics from which the bigger

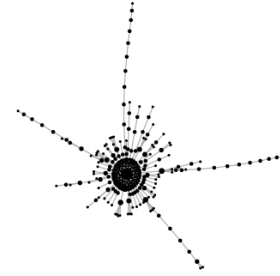


Fig. 1. Representation of discussion threads with post tree graphs. Vertex represent posts and edges represent replies between posts.

graph emerges. On the other hand, to show how structural-temporal neighbourhoods can be used for the detection of different types of conversationalists in online forums or any other type of online discussion that is representable by a tree structure (e.g.: e-mails).

The remaining of the paper is as follows. We first discuss about the convenience of the classic neighbourhood definition in dynamic graphs such as discussion trees. Then we introduce our two definitions of structural-temporal neighbourhoods. To illustrate the kind of neighbourhoods that we obtain in a real data, we analyse the neighbourhood census of a Reddit forum. Finally, we apply both neighbourhood definitions to detect clusters of users that tend to appear in the same neighbourhoods. We close the paper with some suggestions of future research.

II. DISCUSSION TREES

We represent a discussion thread by a tree graph $G = (V, E)$ where V is a set of n vertices representing the posts (also called messages or comments) and E is a set of $n - 1$ directed

edges that say which posts replied to which post. Vertices have three attributes namely order, time, and author. If for two given vertices v_i, v_j there exists an edge $e = (v_i, v_j) \in E$ then v_j is called the *parent* of v_i , denoted as $p(v_i) = v_j$. The *root* of the tree is the only vertex with no parent, and corresponds to the post that starts the discussion. We say that two vertices v_i, v_j are *siblings* if and only if $p(v_i) = p(v_j)$. A vertex v_i is an ancestor of a vertex v_j and v_j is descendant of v_i if and only if v_i is in the path from v_j to the root. A *leaf* is a post with no replies. A *branch* is the path between a leaf and the root. Two vertices v_i and v_j are said to be *neighbours* if either $p(v_i) = v_j$ or $p(v_j) = v_i$ or, in other words, if their distance in the undirected version of G is one.

III. STRUCTURAL NEIGHBOURHOODS

Extending the classic definition of neighbourhood according to which two vertices are neighbours if the distance between them is one, we start by the following definition of *structural neighbourhood*:

Definition 1: Given a tree graph G , the *structural neighbourhood* of radius r of post i , denoted as $\mathcal{N}_i(r)$, is the induced graph composed of all the vertices that are at distance equal or less than r from post i .

In the context of discussion threads, this definition has two limitations. First, the decision on whether to include some post in the neighbourhood is only based on the structural distance, and therefore two posts that are at distance $d \leq r$ are considered neighbours of i regardless of the time when they were written. In conversations, time plays an important role, therefore this is an important limitation. Another consequence of looking only at the structure is that the number of possible neighbourhoods within a radius r is infinite. This poses a problem when trying to categorize conversations since many conversations, while structurally different, can be considered semantically equivalent.

IV. STRUCTURAL-TEMPORAL NEIGHBOURHOODS

As we have discussed, structural neighbourhood falls short in the analysis of conversational structures. Instead, we propose two new definitions of neighbourhood that take into account the order and time in which posts are written.

A. Order-based

Our first definition is based on the order in which posts are attached to the tree.

Definition 2: Given an ordered tree graph G , the *order-based neighbourhood* of radius r of vertex i , denoted as $\mathcal{N}_i^T(r, n)$ is the induced subgraph from its structural neighbourhood composed the n vertices that are closest to i in time and for which there exists a path to i in $\mathcal{N}_i^T(r, n)$.

An example of order-based neighbourhood is given in Figure 2 (left). This definition has two advantages over the *structural neighbourhood*. First, the temporal aspect of the conversation is better taken into account since the neighbourhood only includes posts that are not only near to i in the structure but also in time. Second, the size of the neighbourhood has an

upper bound of $\min(|\mathcal{N}_i(r)|, n)$ thus making the space of possible neighbourhood structures finite. The main limitation of this definition is that the obtained neighbourhoods might be very different with different choices of r and n , and we have no *a priori* criteria to chose the proper parameters other than making them small so that the neighbourhoods capture the local dynamic of the conversation around the post i . We will tackle this issue in the next definition.

B. Time-based

In a first attempt to take time into account, one might set fixed time-based boundaries for the neighbourhood and include only those posts whose timestamp t_j is at distance less than τ from the ego post i , $|t_j - t_i| < \tau$. However, the pace at which posts are added to the conversation may be very different between conversation threads (and also within a thread) and we have no *a priori* criteria for a proper choice of τ . Rather than looking for a fixed time radius, we may look at changes in the pace how new posts are added to the thread. In statistical analysis, a *changepoint* in a sequence x_1, \dots, x_n is a point that comes from a different probability distribution than its precedent values. If applied either to the timestamps of posts in a branch (vertical sequence) or to the timestamps of the replies to the same post (horizontal sequence), we say that a sequence posts with timestamps t_1, \dots, t_n belong to the same (vertical or horizontal) *local dynamic* if there is no changepoint t_i in the sequence such that $1 < i \leq n$. Now we can introduce our second definition of neighbourhood.

Definition 3: Given a tree graph G , the *time-based neighbourhood* of vertex i , denoted as $\mathcal{N}_i^T(r)$, is the maximal subgraph of the structural neighbourhood $\mathcal{N}_i(r)$ where all the vertices belong to the same vertical and horizontal local dynamic than i .

An example of time-based neighbourhood is given in Figure 2 (right). Note that we still have a radius parameter r to guarantee that the resulting structure remains local even if no changepoint is detected near the post i . Algorithm 1 extracts time-based neighbourhoods according to the given definition.

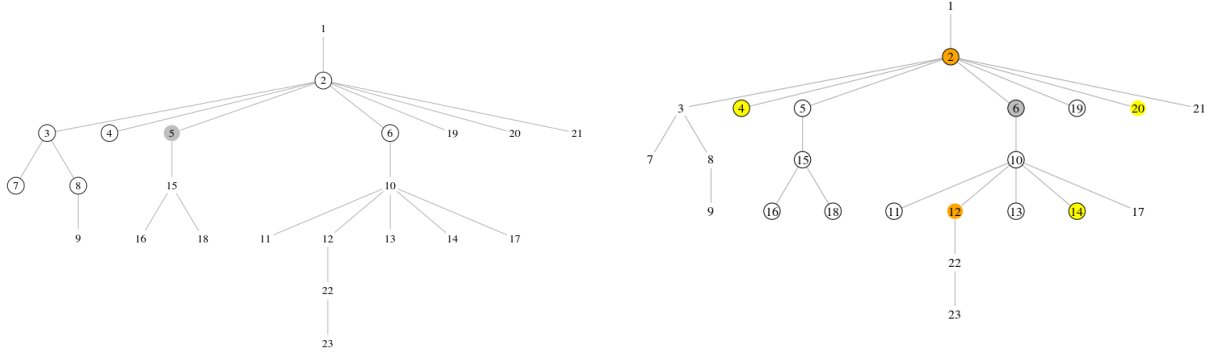


Fig. 2. Illustration of order-based (left) and time-based (right) neighbourhoods. Grey nodes represent the ego post. Nodes with circles are those in the neighbourhood. Numbers indicate the order of the post. The order-based neighbourhood has parameters $r = 3$, $n = 6$. In the time-based neighbourhood, horizontal changepoints (yellow) and vertical changepoints (orange) represent posts that are temporally far from their predecessors (siblings or parents) and therefore set the limits of the neighbourhood.

Algorithm 1 Extraction of time-based neighbourhood

Input: Posts tree g , vertical breakpoints, horizontal breakpoints, ego post ego

Output: Subgraph of g with all vertices in $V(g)$

 Compute structural neighbourhood $\mathcal{N}_i(r)$

 ancestors \leftarrow ancestors(ego) in $\mathcal{N}_i(r)$

 older_siblings \leftarrow older_siblings(ego) in $\mathcal{N}_i(r)$

 dump $\leftarrow \emptyset$

for bp \in vertical breakpoints **do**

if bp \in ancestors **then**

 dump \leftarrow dump \cup ancestors(bp)

else

 dump \leftarrow dump \cup descendants(bp) \cup bp

end if

end for

for bp \in horizontal breakpoints **do**

if bp \in (older_siblings \cup ancestors) **then**

 dump \leftarrow dump \cup older_siblings(bp)

else

 dump \leftarrow dump \cup younger_siblings(bp) \cup bp

end if

end for

$\mathcal{N}_i^{(t)}(r) \leftarrow$ delete(dump_posts) from $\mathcal{N}_i(r)$

C. Neighbourhood colouring

Even if the structure contains some important information about the type of conversation, there is still some ambiguity left, and two similar neighbourhoods can represent very different types of conversation. We can easily reduce this ambiguity by assigning colors, or labels, to vertices, identifying some relevant property of the post.

In particular, we assign the red color to the ego post, the orange colour to all other posts written by the author of the ego (to identify re-entries) a white color to the root

(if included in the neighbourhood), the yellow color to all posts written by the same author who wrote the parent of ego (to identify, for instances, debates between the ego author and the other author), the grey color to all posts written before the ego and the black color to all posts written after the ego. Since some posts might correspond to several colors, we perform a sequential assignment of colors given by grey, black, orange, ego, yellow, white.

D. Neighbourhood pruning

Although the two previous definitions reduce the neighbourhood space considerably, we can still find structures whose only difference is that one of them has more leafs hanging from some of the nodes. Thus, we prune every neighbourhood by leaving a maximum of two consecutive siblings of the same color and where neither of them has any children.

Figure 3 shows some real examples, from our dataset, of time-based neighbourhoods after colouring and pruning.

V. APPLICATION TO REDDIT FORUMS

We applied our two definitions of structural-temporal neighbourhood (order-based and time-based) to a forum dataset. Our dataset consists of all posts from March 2014 to May 2015 of the Podemos forum in Reddit¹. The Podemos forum was conceived in March 2014 as a tool for internal democracy, and forum members used it to debate ideological and organizational principles that were later formalized in their first party congress hold in Madrid the October 18th and 19th 2014. Nowadays, its members use it mainly to share and discuss about political news. Our dataset contains 75,000 posts spread over 4,262 threads and written by 9,160 users between the April 25th and October 20th 2014.

¹<https://www.reddit.com/r/podemos>

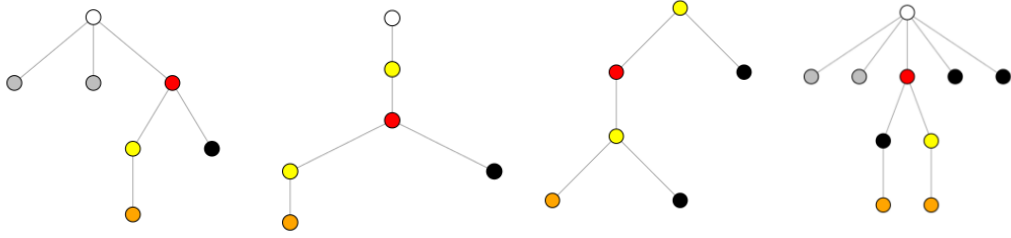


Fig. 3. Illustration of different real structures using different colors for the root (white), the ego (red) and other posts written by the same author (orange) the ego parent and other posts written by the same author (yellow), posts previous to the ego (grey) and posts posterior to the ego.

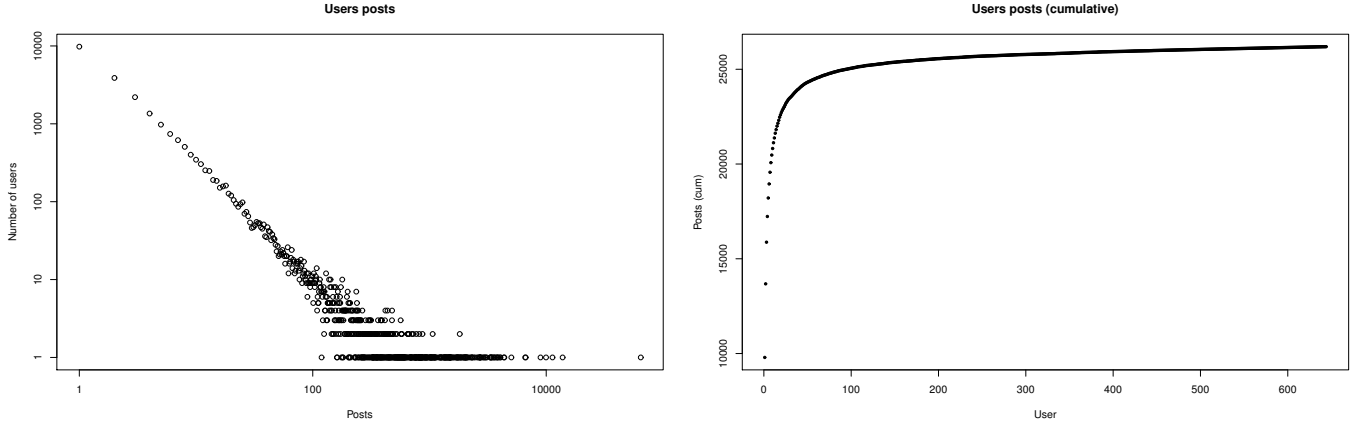


Fig. 4. Log-log and cumulative distributions and of posts per user.

A. Neighbourhood census

Given a set of local structures and a graph under study, a census is a counting of how many times does every structure appear in the graph. In the context of online discussion forums, triad census in the graph of users interactions are often used (in these graphs, an edge exists from user u to user v if u wrote a post reply to some of the posts written by v in some thread.) [6], [7].

We analysed the order-based and time-based neighbourhood census of our dataset. A counter for every neighbourhood structure is created the first time it is detected. For every post, we extract its neighbourhood (order-based or time-based) colour and prune it as explained above (Section IV-C and IV-D). Then we check whether it is isomorphic to some of the already seen neighbourhoods. If this is the case, then we increment the counter for the neighbourhood previously seen. Otherwise we create a new entry for the current neighbourhood.

Figure 5 and 6 show the most frequent order-based and time-based neighbourhoods in the Podemos dataset. In Figure X we compare the census of the most frequent neighbourhoods in the Podemos and the Game of Thrones forums.

B. Conversation-based clustering of users

In this section, we cluster users based on the neighbourhoods of their posts. Our intuition is that some users like

participating in some kind of discussion rather than other. Certainly, most part of the information necessary to understand the nature of a discussion is in its textual content. However, due to the huge diversity of topics, vocabulary, and the difficulty of current algorithms to capture the language subtleties such as humour, irony, or context, we turn our attention towards the structure of the discussions, which can also contain some information. We work under the hypothesis that the structural-neighbourhoods in which a user post is embedded reflect the kind of conversation in that part of the thread.

We limit our study to a set of 100 *active users* who wrote more than 100 posts. For those users, we create an initial feature matrix $U \times N$ where U is the number of users and N is the number of neighbourhoods in the census, and where the position (u, n) is a counter of the number of times that a post written by user u has a neighbourhood isomorphic to n . We drop those feature columns that are zero for every user (these are neighbourhoods seen only around posts of users with low-level activity). To make the feature vector of a user independent on the number of posts, we transform the counts into percentages. And since some features have much higher percentages than others for most users, we scale and normalize the matrix so that every feature has mean 0 and variance 1. To avoid non-significant scores, we remove also the feature column corresponding to neighbourhoods that have a frequency less than 50 among the active users.

We use k-means to find the clusters, though one can use any

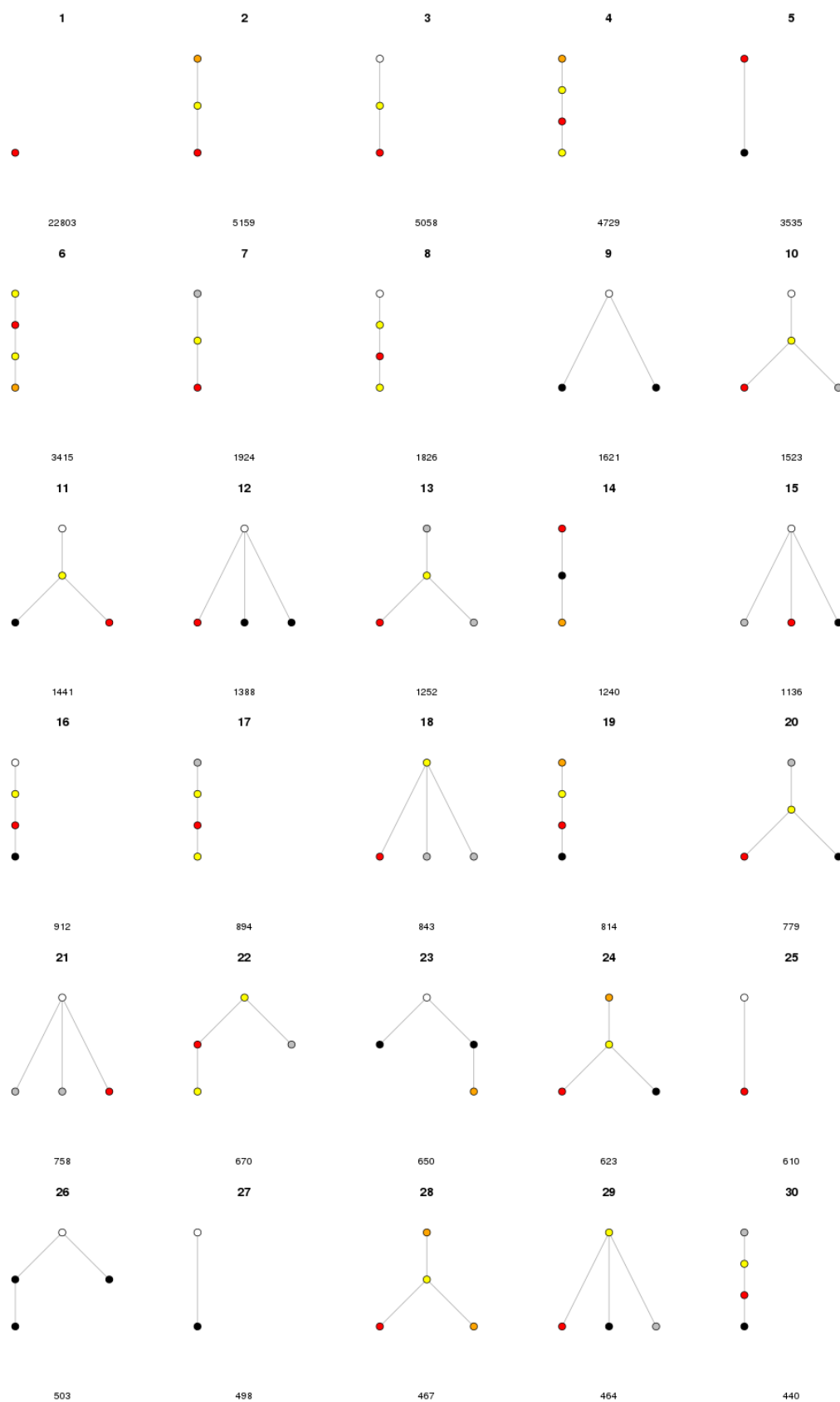


Fig. 5. Most frequent order-based neighbourhoods with $r = 2$ and $n = 4$

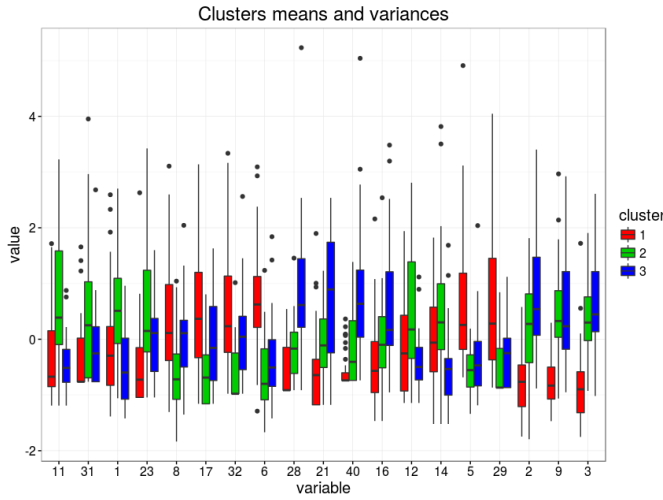


Fig. 8. Whiskers plots in the most relevant features (time-based)

other clustering method. Since the plot over the Within-Cluster Sum of Squares for $k = 1, \dots, 20$ did not show any clear elbow we chose $k = 3$ clusters so that clusters are easier to interpret. We did the clustering over a feature matrix with order-based neighbourhoods and another feature matrix with time-based neighbourhoods. Figure 7 and shows the PCA projections of the users coloured by their assigned cluster.

VI. CONCLUSIONS

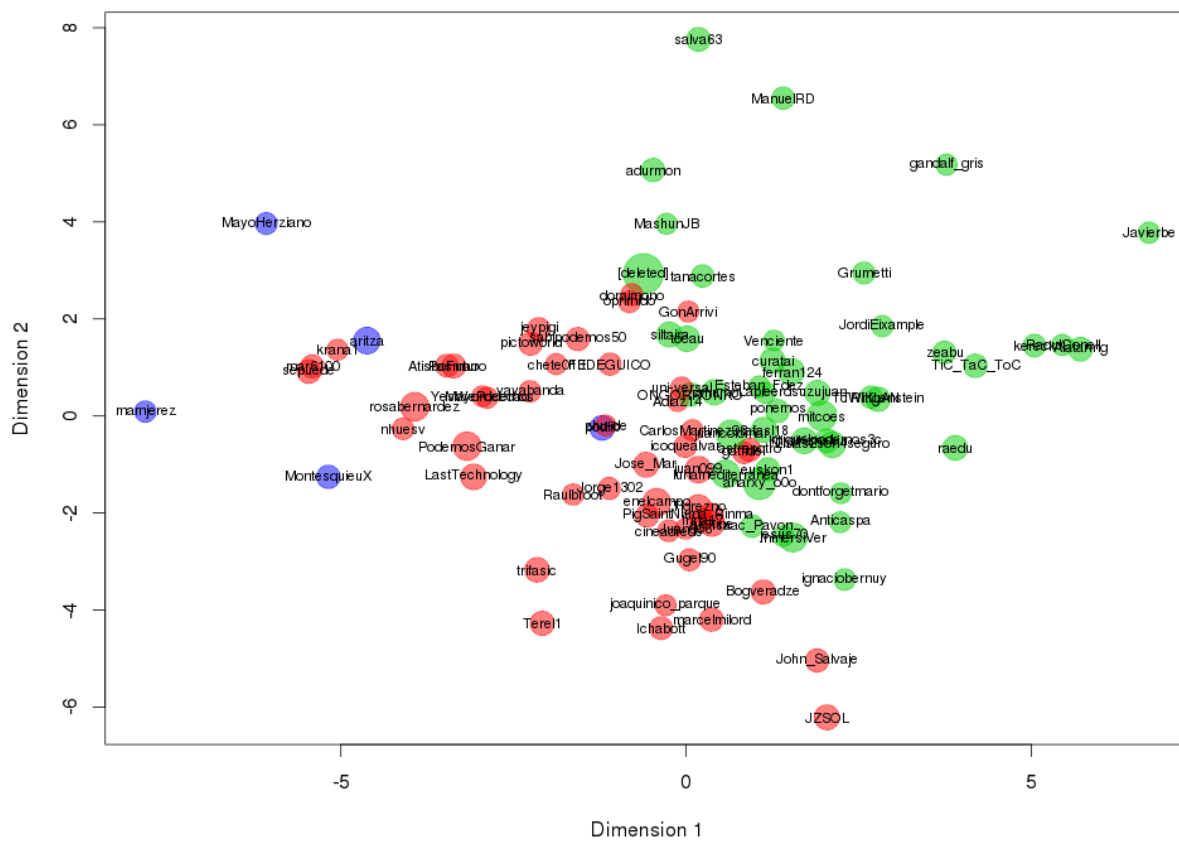
We presented a method to characterise conversations of users in online threads. Due to the tree nature of online threads, traditional patterns such as triads are not able to capture much of relevant dynamics of a conversation. Our defined order-based and temporal-based neighbourhoods are able to capture a very rich variety of structures. We used this neighbourhoods to characterise users in terms of the structure of the conversations they participate in and showed that, indeed, there are different types of structural conversationalists

The concept of structural-temporal neighbourhood opens the door to some interesting paths of research. One might wonder whether other pruning are more pertinent than the proposed here. Also, even after pruning some neighbourhoods suggest the same type of conversation, so a manual merge might be convenient. Maybe sociologists can help to merge the found neighbourhoods in a meaningful way.

REFERENCES

- [1] A. McCallum, X. Wang, and A. Corrada-Emmanuel, "Topic and role discovery in social networks with experiments on enron and academic email," *Journal of Artificial Intelligence Research*, vol. 30, no. 1, pp. 249–272, 2007. [Online]. Available: <http://www.aaai.org/Papers/JAIR/Vol30/JAIR-3007.pdf>
- [2] R. Bhatt and K. Barman, "Global Dynamics of Online Group Conversations," in *Proceedings of the Sixth International AAAI Conference on Weblogs and Social Media*, 2012.
- [3] R. Dorat, M. Latapy, B. Conein, and N. Auray, "Multi-level analysis of an interaction network between individuals in a mailing-list," *Annales des télécommunications*, vol. 62, no. 3-4, pp. 325–349, 2007. [Online]. Available: <http://link.springer.com/article/10.1007/BF03253264>
- [4] S. Whittaker, L. Terveen, W. Hill, and L. Cherny, "The dynamics of mass interaction," in *Proceedings of the 1998 ACM conference on Computer supported cooperative work*, 1998, pp. 257–264.
- [5] N. Gaumont, T. Viard, R. Fournier-S'niehotta, Q. Wang, and M. Latapy, "Analysis of the temporal and structural features of threads in a mailing-list," in *Workshop on Complex Networks CompleNet*, 2016. [Online]. Available: <http://arxiv.org/pdf/1512.05002v1.pdf>
- [6] L. A. Adamic, J. Zhang, E. Bakshy, and M. S. Ackerman, "Knowledge sharing and yahoo answers," in *Proceeding of the 17th international conference on World Wide Web - WWW '08*. New York, New York, USA: ACM Press, 2008, p. 665. [Online]. Available: <http://dl.acm.org/citation.cfm?id=1367587> <http://portal.acm.org/citation.cfm?doid=1367497.1367587>
- [7] A. Lumbreras, J. Lanagan, J. Velcin, and B. Jouve, "Analyse des rôles dans les communautés virtuelles : définitions et premières expérimentations sur IMDb," in *Modèles et Analyses Réseau : Approches Mathématiques et Informatiques (MARAMI)*, 2013, pp. 1–12.

Individual factor map (PCA)



Individual factor map (PCA)

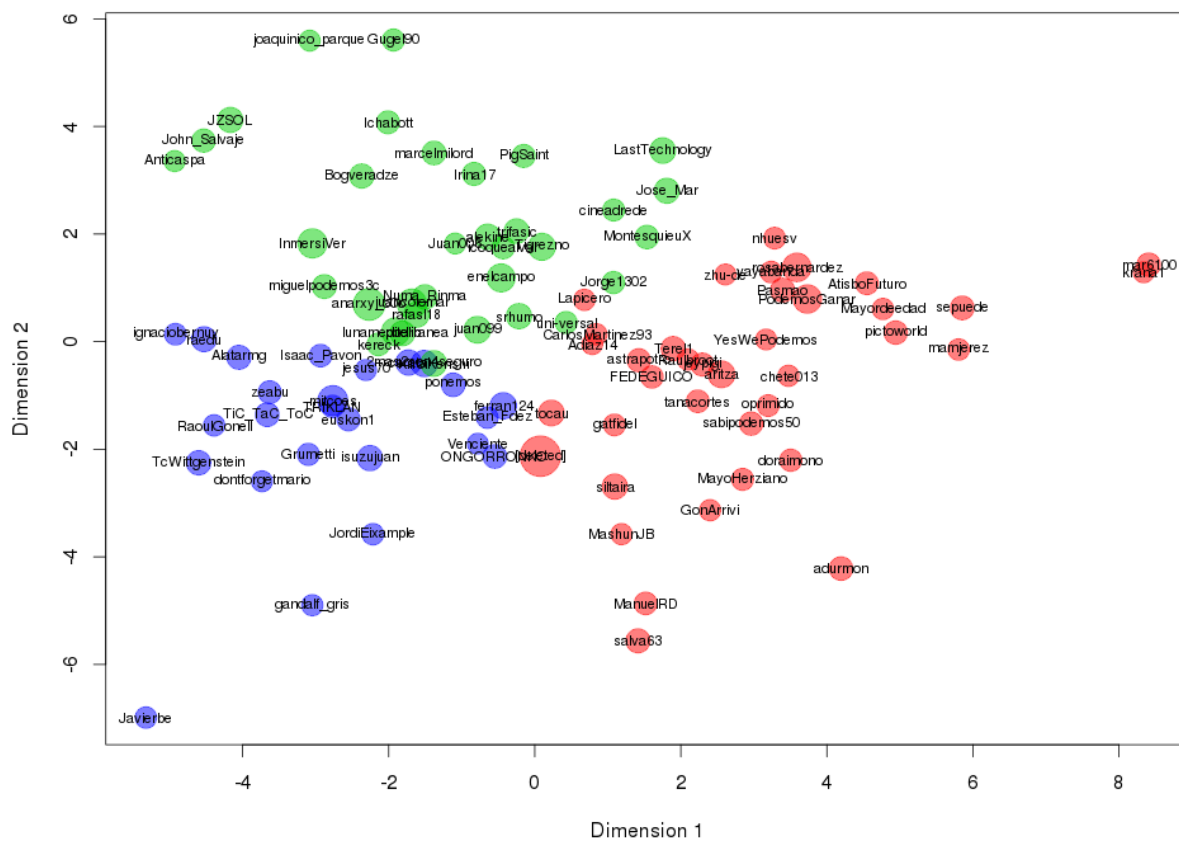


Fig. 7. PCA projections of the order-based (above) and time-based (below) neighbourhood features