

Modeling User Reports in Crowdmaps as a Complex Network

Carlos Caminha

Universidade de Fortaleza
Av. Washington Soares 1321, Fortaleza, CE, Brazil
+558534773287

carlos.o.c.neto@gmail.com

Vasco Furtado

Universidade de Fortaleza
Av. Washington Soares 1321, Fortaleza, CE, Brazil
+558534773287

vasco@unifor.br

ABSTRACT

Collaboration with content sharing via digital maps is a type of application that is characteristic of the context of the social web. In these applications, the collaborative map works as a blackboard for accommodating descriptions of events to be shared, typically with members of a social network. A malicious activity that is difficult to detect in this interactive context is the generation of a false trend on the map as the result of several false reports by more than one person. In this article, we outline the general lines of our investigation into how modeling in complex networks of events registered on a collaborative map can help identify regularities, and therefore show deviations arising from malicious activity. In particular, we will focus on the context of Public Safety, with an analysis of the distribution of crimes reported by the population in a geographic area. The idea here is to model a network comprised of users who reported events and the locations where such events were reported (e.g.: a census tract).

Keywords

Complex networks, security, crowdmaps.

1. INTRODUCTION

Motivated by the huge success of Wikipedia, wiki applications have not been restricted to crowdsourcing via text sharing. On the contrary, there has recently been an explosion of interest in using the web to create, assemble, and disseminate geographic information provided voluntarily by individuals. Crowd mapping, combining the aggregation of a Geographic Information System and crowd-generated content, flourishes daily on the Web [6], [9]. Sites such as Wikimapia (<http://www.wikimapia.com>), Click2fix (<http://www.click2fix.co.sa>), Crowdmap (www.crowdmap.com), and OpenStreetMap (<http://www.openstreetmap.org>) are empowering citizens to create a global patchwork of geographic information, while Google Earth and other virtual globes are encouraging volunteers to develop

interesting applications using their own data. In crowd map applications, the digital map works as a blackboard for accommodating stories told by people about events they want to share with others typically participating in their social networks.

Our research is inserted into this context. We have developed a platform to create and host crowdmaps called WikiMapps (www.wikimapps.com) [2]. Several maps have been created in WikiMapps for different domains. Among them, one with large popularity is WikiCrimes (www.wikicrimes.org) [5], which provides a common area of interaction among people so that they can report and monitor the locations where crimes are occurring. WikiCrimes allows users to access and to register criminal events on the computer directly in a specific geographic location represented by a map. Alarms that indicate the most risky places and heat maps are example of services produced by the website for people in general.

In this article we concentrate on a fundamental aspect within the WikiCrimes context and the crowd mapping in general: the identification of trends that show to be malicious activities coming from hoaxes or false reports. An approach for mitigating this problem is to attribute a credibility score to the information based on a model of trust and reputation of the users [5], [12]. To take into account the social network of the users for community detection and consequently outliers of these is another possibility [11].

However in crowdmaps there is a need to keep in balance the trade-off between diminishing the constraints imposed to the users with the intention to increase the number of participants in the system, and the rigid control that can be imposed to avoid unwanted behavior, such as the reporting of false information. For this reason, little information about the users is available, which ultimately makes the

above-mentioned approaches unfeasible, or makes them produce a very large number of false positives.

Moreover, some malicious activities are not done exclusively by the report of a single user, which could be more easily detected by anomaly detection techniques [1] LOF. False trends that are more difficult to identify are those caused by a group of users (which actually can be done by a single user with several fakes) that report crimes in a certain place with the aim of highlighting it in comparison to others. These malicious actions cannot be captured only through analysis of reports or of users individually; they require an investigation from the perspective of the relationships among users.

This motivated us to consider the exploration of complex networks modeled after the information of the users, the reports, and the locations where the reports were made. This model makes use of patterns identified in previous work [3], [8], showing that the distribution of crimes by census tract follows a power law. It is verified in this context that there are few places that concentrate many crimes, and many places that concentrate few crimes. On the other hand, the literature on collaborative systems has shown that people's participation in collaborative systems such as crowdmaps also has a skewed distribution, which is popularly called the 90-9-1 rule [10]. Many users participate little, and few users participate very actively.

Our approach in this paper is first to characterize the data described in WikiCrimes, which led us to investigate the existence of power law distributions suggested in the literature. Then, we were able to identify new regularities that are evident, particularly with regard to the correlation between users who report crimes in certain places. Starting from a bipartite network model in which the vertices are individuals and census tracts, we projected a monopartite network of users in which the edges indicate the strength of connection between them. This connection strength indicates the degree of co-relatedness of the reports of crime made by these two users in a particular place.

Based on this modeling and on information obtained by the characterization of the data such as the distribution of crime per census tract and the distribution of reports from users, we were able to find a certain regularity within the context of WikiCrimes. This regularity refers to the fact that hubs have a high geographic coverage (i.e. they report

crimes in the majority of census tracts) and therefore demonstrate a well-defined behavior with regard to their connections with non-hub users. By characterizing this, we were able to observe that the relationships of non-hub users among themselves are typically no stronger than the relationship between such non-hub users and the hubs. If this happens, the possibility of malicious activity becomes abundantly clear. In this article we overview our method.

2. RELATED WORK

Wikis, in general, are based on the concept of radical trust; i.e., it is believed that individual participation, for the most part, includes correct information. Nevertheless, the identification of attempted fraud or vandalism is necessary. The challenge imposed on WikiCrimes and collaborative maps in general is to assure the credibility of the information recorded on the map, and requires the study of different approaches.

One approach to minimizing the problem is to assign a value to the credibility of the information based on a model of user reputation and trust. The analysis of the users' social network, for example, was one of the ways we used to achieve this [5] and [12]. Reputation models, however, lack the level of granularity to capture malicious activities such as generation of a false trend that can come about with an excess of false reports made by various people. Identifying evidence of these problems through data mining is another recommended approach.

Based on this, we developed an algorithm [4] that – based on reported events – tries to identify patterns that indicate excesses or abuses coming from an individual or group of individuals. The basic idea is to identify the existence of communities on the social network and verify if there is one such community that dominates the reporting of events (reported for a hot spot, in particular). In order to do this, we used algorithms for identifying communities developed in the context of social network analysis. In addition to considering the structure of the community in terms of connection density, it was necessary to consider the participation of users (represented in the nodes) in the formation of hot spots. Since the formation of hot spots varies with the zoom level, the authors proposed identifying communities for each one of the social networks related to the hot spots at each zoom level.

Although the first results obtained with this approach were satisfactory, they soon showed that a high number of false positives could occur as a greater

participation in the system starts to occur. Groups of users can report information and generate false trends in a region without forming a community among themselves. Assuming that most of the reports are true, one must try to find this type of anomaly, if any.

In anomaly detection, the goal is to find objects having behavior that is very different or extraneous in relation to others. In the context of WikiCrimes, these objects can be users, whereby the intent is to extract characteristics of normal behavior thereof, then identify extraneous elements.

The task of detecting anomalies brings many challenges, as it may be necessary to use several or just one attribute to detect one of these elements. These objects can also be anomalous to a certain degree; therefore it could be interesting not to use only a binary definition whereby an element is either extraneous or not. Elements with low degree of anomaly can be extremely difficult to identify. Also, anomalies may not always occur from an overall perspective, since an element may appear normal in relation to the entire network, but on the other hand, it can be extraneous in its vicinity or region in which it is located.

Regardless of the technique chosen to detect extraneous elements, such technique will always require the selection of variables, attributes or classes as input for the use thereof. One of the contributions of this study is to propose a way to represent the collaboration of crime reports in places such as a complex network, by modeling the relations between users in such a way as to make it possible to separate malevolent users from users with desirable behavior within the system.

3. MODELING OF COLLABORATION AS A COMPLEX NETWORK

3.1 Representation of reports of crime in census tracts

On collaborative maps, bits of information are marked by users in different geographical regions. Specifically in WikiCrimes, users report occurrences of various types of crimes anywhere in the world on a digital map. Typically the map has a kernel layer [7] indicating the density of crimes. Redder colors indicate a high concentration of crime (also called hot spots).

Hence, our complex network model is based on

information from users, reports of crimes made by such users, the locations where the reports refer to, represented here by census tracts. This model is based on a bipartite graph and its projection is specified below.

The directed bipartite graph $G^b(U, S, E^b)$ has – as vertices – the users $u (\in U)$ and the census tracts $s (\in S)$. An edge, e^b , represents the fact that there was a report by a user u in a tract s . The weight of the edge $e^b (\in E^b)$ is obtained from the number of crime reports made in s .

In order to have a representation that indicates the strength of the connection between the users, we projected the bipartite graph onto a monopartite graph in which vertices are the users. Thus, users who reported crimes in the same tract will have an edge that joins them, in which the weight of this edge is the number of crimes they reported in common in that location. In other words, if users report crimes in more than one tract, the number of crimes reported by them in common in these tracts is added to the weight of the edge. Formally, the monopartite graph $G(U, E)$ has on its vertices the users $u (\in U)$ and the edges $e (\in E)$. The weight of an edge e , $w(e_{uu'})$, between two users u and $u' (\in U)$ is calculated based on the weight of the edges of u with $s (\in S)$ and of u' with $s (\in S)$ in the bipartite graph G^b . More specifically, $w(e_{uu'})$ is the sum of the minimum number of crimes reported by u and u' , in all tracts, according to the formula below:

$$w(e_{uu'}) = \sum_{i=1}^n \left(\min \left(w(e_{us_i}^b), w(e_{u's_i}^b) \right) \right)$$

where n is the number of census tracts, s_i , in which both u and u' reported crimes.

The strength of the relationship between users, represented by the weight of the edge $w(e)$, indicates that the higher the weight, the more those users reported crimes in the same locations.

3.2 Characterization of WikiCrimes data

Modeling user interaction as a complex network allows one to extract the main properties of such network to better understand how users relate to one another. This analysis indicates that the process of forming this network does not seem random. On the contrary, the characteristics extracted from this network show a preferential connection process of the nodes, where new vertices inserted tend to connect

with the hubs of the network. Some of the properties extracted were the following:

3.2.1 Graph Hub

The network hub has degree 352, i.e., this vertex has nearly 20% of all the network's edges, which total 1768. It is also interesting to see that it connects with 352 vertices, and this value represents more than 90% of the total number of nodes on the network. This stems from the fact that the hub reported crimes in almost all the census tracts (more precisely: 82% of the tracts).

3.2.2 Degree Centrality

The degree centrality of this network is 0.891; this value shows how easy it is to break it. We conducted an experiment where preferentially removing the higher-degree vertices, we were able to "break" more than 90% of the network by removing only 10% of the nodes.

3.2.3 Network Density

The network has a low density, with only 2.35% of the maximum density that a network of 388 nodes can attain. This value shows the great distance between the degree value of the graph's hubs and the degree value of the vertices with fewer adjacencies on the network.

3.2.4 Network Diameter

The shortest path between all vertices is, at most, four leaps. The hubs seem to play an important role in maintaining this property at such a low value; the fact that they report crimes in many places makes the shortest paths between users that report in few places remain very small.

3.2.5 Clustering Coefficient

The clustering coefficient of the network is 0.93. This high value is due to the projection we made from the bipartite graph to a monopartite graph, with complete sub-networks, because all users who report crimes in the same location have edges connecting them all together.

We conducted an analysis with a particular focus of attention on the distribution of the degrees and distribution of the ranking of degrees. In this analysis, we saw that the number of vertices with degree 2, 3, 4, 5 and 6, together represent 61% of the total number of network nodes. This value shows that many vertices in such network have a low degree and few have a very high degree, which explains the low density in item 3.1.3 and the diameter of only four leaps in item 3.1.4. The degree distribution has a

slope -1.655. We ordered the vertices by the degree value and plotted the ranking position of each vertex on the "x" axis and the degree thereof on the "y" axis. Figure 1 shows this graph. It indicates a slope of -0.839 with R^2 of -0.950. These figures show that in spite of hubs with many connections in relation to the lower-degree nodes, the increase in the value of these degrees takes place in a way that is milder than expected, therefore has a relatively low slope. We also analyzed the distribution of weights of the edges. This allowed us to discover that this is also a power law, which in this case assumes a slope of -1.56. Figure 2 illustrates the curve of this distribution, where we plotted the frequency at which each edge value is repeated on the "y" axis and the repetition frequency ranking on the "x" axis. The values of the distributions in this figure are in log scale.

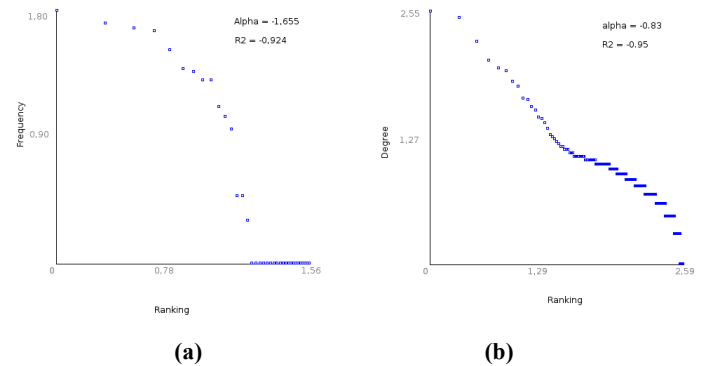


Figure 1. Graphs indicating the degree distribution (a) and frequency ranking (b)

This analysis also led us to identify something that proved very important to be considered in identifying malicious activity: we found that for each frequency value plotted, there was at least one of the hubs that had that degree represented. A "hub" is understood as a vertex belonging to the set of vertices with the highest degree in the network, where the sum of the edges of the members of this set adds up to 80% of the edges of the network. This proved to be particularly relevant information, as we reinforce below.

4. IDENTIFICATION OF MALICIOUS ACTIVITY

The properties extracted from this network make it seem like a *small world*, with low diameter and high clustering coefficient. The hubs are very connected, the result of comprehensive activity that exists in the reporting of events in various census tracts.

The explanation for this fact is that in collaborative systems of crime reports, where official data and data coming directly from the population are mixed, the hubs are typically entities or people with good reputations, such as government agencies, and that possess information that is mapped on a wide geographical area. In WikiCrimes, when analyzing data from a large Brazilian city, we saw that the hubs are the so-called “certifier entities” [5], project partners that hold a large volume of information on crimes, such as insurance brokers, police officers, specialized media, etc. Because they have large quantities of data, they make their reports in various census tracts in the city.

These hubs have an essential role in the behavior pattern of users and their reports of crimes, and seemed to be the key for detection of activities that may indicate fraud.

It is noteworthy that in the monopartite network of users, the weight of the edge determines the degree to which two users, connected by that edge, report events in the same places. Analyzing the relationship of these users, we noted that, in all of them, the edge with the highest weight is one that refers to its direct connection to a hub. In the case of WikiCrimes, we’re talking about a set of only 14 users. Here we saw a type of regularity that is evident of this type of scenario, i.e., non-hub users report crimes in areas where hubs also do. In other words, the coincidence of reports in a given area is much more likely to occur with reports made by hubs rather than with reports made by other, non-hub users.

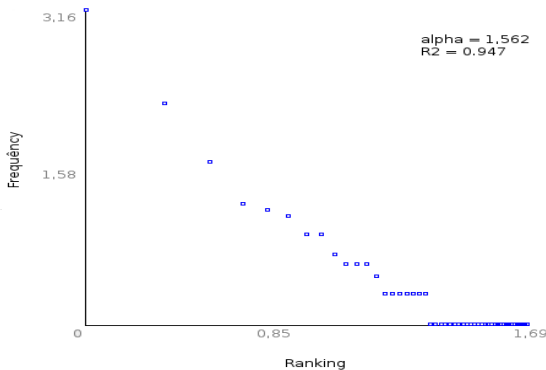


Figure 3. Plot of ranking by frequency at which the edge weight values are repeated

This relationship pattern between non-hub users and hubs proved to be worthy of verification, because the absence of this pattern in the behavior may indicate suspicious relationship between users, whereby

people who do not have a volume of reporting sufficient to generate a false trend map, start to have this power by adding their entries to those of other users at that location.

Formally, verification of this pattern can be described as follows. Let U_h ($U_h \subseteq U$) be the set of users u_h with a high degree, ($d(u_h) > \phi$ where ϕ is a system parameter) and a triangle T formed by a hub ($u_h \in U_h$) and by two users u_1 and u_2 ($\in U$). The correlation factor $\rho(u_1)$ of a vertex of the triangle is a ratio of the weight of the edge $w(e_{u_1u_2})$ of this vertex with another user by the weight of the edge with the hub, $w(e_{u_1u_h})$, as expressed in the following formula:

$$\rho(u_1) = w(e_{u_1u_2})/w(e_{u_1u_h})$$

Based on the correlation factor of users in these triangles, one can see the limit for a given vertex (user) to be considered anomalous in its activity of reporting events on maps. We verified that, for the WikiCrimes network, it is normal to expect values of $\rho \leq 1$.

An analysis on the WikiCrimes base did not lead to the identification of any malicious activity by a group of users. Through simulations and analysis of the behavior of the correlation factor, we found that our approach is a good strategy to create alerts that malicious activity may be occurring.

5. CONCLUSION

Our scientific research has sought to represent the participation in crowdmaps through a complex network. By doing so, we can better understand the patterns that form based on the relationships between users who report events on the maps. Thus, we were able to formally measure the relationships between users and identify patterns of the behavior thereof within the network. The hubs of this network are the key to detecting anomalies. This stems from the fact that the hubs are usually entities with a high reputation and very often refer to government agencies. These are users who participate actively in the reporting of events and do so in several places, exposing a clear pattern of relationship between them and other users.

In particular, by following this approach, we were able to develop a method for detecting a type of malicious activity that is extremely difficult to identify in the context of crowdmaps. We were able to identify that relationships between non-hub users among themselves are typically no stronger than

relationships between non-hub users and the hubs. When this occurs, the possibility of malicious activity becomes strongly evident. We formalized how this can be evidenced by means of a measurement of correlation between non-hub users and hubs.

Our attention is now on developing a method to generate networks with the characteristics that we identified as existing on crowdmaps, i.e., users report occurrences following a power law; places receive reports of occurrences following a power law; and the quantity of places for which users report occurrences also follows a power law.

6. ACKNOWLEDGEMENT

The second author was partially funded by CNPq grant: 55977/2010 and 304347/2011.

7. REFERENCES

- [1] Breunig, M. M.; Kriegel, H. -P.; Ng, R. T.; Sander, J. (2000). "LOF: Identifying Density-based Local Outliers". *ACM SIGMOD Record* **29**: 93. doi:10.1145/335191.335388.
- [2] Caminha, C. Furtado, V. Ayres, L. Vasconcelos, E.: Uma ferramenta de autoria para criação de mapas colaborativos para aplicações em egov 2.0, Anais do WCGE 2010, SBC, 2010.
- [3] Cançado, T.: Alocação e despacho de recursos para o combate à criminalidade. Dissertação de mestrado, Departamento de Ciência da Computação, UFMG, 2005.
- [4] Furtado, V. Assunção, T. de Oliveira, M., Belchior, M. D'Orleans, J.: A Method for identifying Malicious Activity in Collaborative Systems with Maps. Algorithms for Social Network Mining (ASONAM), Athens, 2009.
- [5] Furtado, V. Ayres, L. de Oliveira, M. Vasconcelos, E., Caminha, C., D'Orleans, J.: Collective Intelligence in Law Enforcement: The WikiCrimes System, Information Science, 180, 2010.
- [6] Mac Gillavry. Collaborative Mapping and GIS: An Alternative Geographic Information Framework. In: Balram, S. & S. Dragicevic (eds.) Collaborative Geographic Information Systems, pp. 103-119. London: Idea Group Publishing, 2006.
- [7] McLafferty, S.: Identification, development and implementation of innovative crime mapping techniques and spatial analysis. Washington, D.C: U.S. Department of Justice, 2000, p.27.
- [8] Melo, A.: Um modelo multiagente de simulação criminal bio-inspirado. Dissertação de mestrado, Mestrado em Informática Aplicada, UNIFOR, 2008.
- [9] Rouse, J., Bergeron, S.J., Harris, T.M., "Participating in the geospatial web: collaborative mapping, social networks and participatory GIS," in Scharl, A., Tochtermann, K. (Eds.), The Geospatial Web: How Geobrowsers, Social Software and the Web 2.0 are Shaping the Network Society, Springer, New York, NY, pp. 153-8, 2007.
- [10] Whittaker, S. Terveen, L. Hill, W. and Cherny, L. (1998): "The dynamics of mass interaction," *Proceedings of CSCW 98, the ACM Conference on Computer-Supported Cooperative Work* (Seattle, WA, November 14-18, 1998), pp. 257-264.
- [11] Wu, F. and Huberman, B. Finding communities in linear time: a physics approach. The European Physics Journal B, 2004.
- [12] Wu, B., Goel, V., Davison, B.D.: Propagating trust and distrust to demote web spam, in: Proceedings of Models of Trust for the Web Workshop – MTW, Edinburgh, Scotland, 2006.