# Analyzing the Effectiveness of Graph Metrics for Anomaly Detection in Online Social Networks

Reza Hassanzadeh, Richi Nayak, Douglas Stebila

School of Electrical Engineering and Computer Science, Science and Engineering Faculty,
Queensland University of Technology, Brisbane, Australia
{r.hassanzadeh,r.nayak,Stebila}@qut.edu.au

**Abstract.** Online social networks can be modelled as graphs; in this paper, we analyze the use of graph metrics for identifying users with anomalous relationships to other users. A framework is proposed for analyzing the effectiveness of various graph theoretic properties such as the number of neighbouring nodes and edges, betweenness centrality, and community cohesiveness in detecting anomalous users. Experimental results on real-world data collected from online social networks show that the majority of users typically have friends who are friends themselves, whereas anomalous users' graphs typically do not follow this common rule. Empirical analysis also shows that the relationship between average betweenness centrality and edges identifies anomalies more accurately than other approaches.

**Keywords:** Anomaly detection, Graph mining, Data mining, Online Social Networks.

## 1    Introduction

Online social networks are being used in various domains such as business, education, telemarketing and many others. With increasing use of social networks comes increasing prevalence of illegal activities using social networks [1]. It is critical that methods of anomaly detection in social networks are developed to coincide with developments in usage of social networks.

An online social network can be modelled as graph [2] in which the nodes represent people and the edges represent the links between nodes using a range of relationships such as friendship, affiliation, family and many others. In this paper we propose the use of various graph properties for differentiating people's online behaviour by their usage patterns. If the usage pattern of a user follows common patterns, we describe the usage as *normal*, otherwise the usage is an *outlier* or *anomalous*. Looking at the relationships of users can reveal meaningful patterns: users can hide their identity by supplying false information but they cannot hide certain types of metadata, such as the links that they have established with other users.

We use *local graph properties* to extract common rules. Local metrics refer to a single node (e*go*), its 1-level neighbourhood (an *egonet*) and 2-level neighbourhood (a *super-egonet*). These undeniable relationships can help in spotting behaviours that are abnormal.

In particular, we propose the use of *betweenness centrality* and *average betweenness centrality* of a user's egonet, and the *community cohesiveness* of the user's super-egonet as potential measures for identifying anomalies based on the structure of users' links. Additionally, we give a framework for evaluating the effectiveness of various combinations of properties for identifying anomalous nodes in unlabelled datasets.

We evaluate the proposed methods with existing data collected from three online social networks (Facebook, Orkut, and Flickr). Results show that the majority of users follow the "friends of friends are often friends" pattern and a very few users follow either the "cliques or near-cliques" pattern (all the neighbours connected) or the "stars or near-star" pattern (mostly disconnected). Previous works [1, 3-5] have established that these two types of patterns can be connected to abnormalities in the graph, particularly in online social networks. Several graph theoretic metrics, in particular average betweenness centrality give better accuracy in detecting anomalies than existing approaches.

## 1.1 Related Work

Limited work has been done on applying anomaly detection techniques to online social networks until recently [4, 6]. Recent work can be divided into two categories: *behaviour-based techniques* that consider the dynamic usage behaviour of users; and *structure-based techniques* [1, 3-5, 7] that consider the static structure of the graph. Behaviour-based techniques concentrate on mining users' usage patterns. Although they can help to spot anomalies, they are very technology-dependent. Akoglu et al [4] designed a structure-based approach entitled the OddBall algorithm for analyzing social network graphs. OddBall is based on the power law relation between number of nodes and number of edges and a density-based outlier detection technique to calculate a final anomaly score. However, using only a power law relation is prone to miss some outliers especially for egonets with a high number of nodes and edges. In this paper, we propose an algorithm and a framework for detecting anomalies in an unlabelled social network's dataset based on betweenness centrality.

In traditional data mining, a common method of detecting outlier is identification of *clusters*. Similarly, within the modelled graph, a *community* can be defined as a group of nodes which share common properties. Detecting communities can give us useful information to find if there are any similarities or common interests between the friends of suspected users. Graph-based community detection techniques have been investigated in the literature [8, 9]. Existing algorithms try to find parts of the network that are better connected internally. We propose an alternative method based on the number of external links between two users' egonets.

# 2    The Proposed Framework

We propose a framework that introduces semi-supervised graph-based anomaly detection with the use of a scoring method to report anomalies. It aims to find the common behaviour that is followed by the majority of nodes. It computes graph metrics of a user's egonet and then examines relationships between these properties. The common patterns are then used in distinguishing users that may be anomalous. Our proposed analysis method consists of the following steps, which will be explained in detail in the rest of this section:

**Step 1: Compute graph metrics**
   Metrics computed include: N: number of nodes in a user's egonet; E: number of edges in a user's egonet; ABC: the average betweenness centrality of all nodes in a user's egonet; and Com: the community cohesiveness of the user's super-egonet.

**Step 2: Compute fitting curve**
   For the relationships between N vs. E, ABC vs. E, and N vs. Com, the fitting curve will be computed. The fit may be linear or power law [10].

**Step 3: Compute outlier score**
   For each relationship, an anomaly score function, which is based on distance from the fitting line, is determined.

**Step 4: Label for evaluation**
   A labelled subset of nodes is obtained.

**Step 5: Find threshold**
   Using the scoring function from step 3, a threshold that minimizes the number of false negatives and false positives rate is determined for the labelled subset of data.

## 2.1    Step 1: Compute Graph Metrics

A graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ consists of a set $\mathcal{V}$ of vertices (nodes or users) and a set $\mathcal{E}$ of edges (links between two users). Given the graph $\mathcal{G}$, an ego $i$ is a user (or node) and $egonet_i = \{i, i_1, i_2, i_3, \dots, i_n\}$ consists of the user's neighbours $i_1, \dots, i_n$. A user's *super-egonet* includes the user's egonet and the egonets of all its neighbours: the super-egonet of ego $i$ is $super\text{-}egonet_i = \{egonet_i, egonet_{i_1}, \dots, egonet_{i_n}\}$.

### 2.1.1 Average Betweenness Centrality

The *betweenness centrality* ($\mathcal{BC}$) of a node in a graph is the number of shortest paths between all pairs of nodes within that graph that go through that node.

**Definition 1 (Betweenness centrality).** The *betweenness centrality* of a vertex $i \in \mathcal{V}(\mathcal{G})$ is

$$\mathcal{BC}_i = \sum_{s \neq i \neq d} \psi_i^{sd} / n_{sd} \qquad i, s, d \in \mathcal{V} \tag{1}$$

where $\psi_i^{sd}$ is the number of shortest paths between $s$ and $d$ passing through node $i$ and $n_{sd}$ to be the total number of shortest paths from $s$ to $d$. Brandes' algorithm for computing betweenness centrality runs in time $\mathcal{O}(nm)$ and space $\mathcal{O}(n + m)$,

where $n$ is the number of nodes and $m$ is the number of edges [11]. Each new edge defining a new shortest path will reduce $\mathcal{BC}$ of the central node

We propose the use of the *average betweenness centrality* ($\mathcal{ABC}$) of a node within the node's egonet. Recall that $\mathcal{BC}$ for each node is computed as the number of shortest paths between all pairs of nodes within the egonet that go through that node. Adding edges between nodes in the egonet reduces the betweenness centrality of the ego. Intuitively, an egonet has higher average betweenness centrality when more nodes are involved in shortest paths.

**Definition 2 (Average Betweenness Centrality).** The *average betweenness centrality* of $egonet_i$ is :

$$\sigma_i^{abc} = \frac{\mathfrak{f}(i) + \sum_{j=1}^{n} \mathfrak{f}(i_j)}{n}, \quad where \quad n = |\mathcal{V}^{egonet_i}|, i \in \mathcal{V}(\mathcal{G}) \tag{2}$$

We define $\mathfrak{f}(i_j): \mathcal{V}^{egonet_i} \to \mathbb{R}^{\geq 0}$ as the function that maps each node $i_j$ within $egonet_i$ to its betweenness centrality within its own egonet.

### 2.1.2 Community Detection

People naturally tend to form communities based on their similarity and common interests. This behaviour stands true in online social networks [9]. The information which can be extracted from communities' structure is useful to analyze the behaviour of a user and can lead towards identifying anomalous behaviour. For community detection we examine users' super-egonets, which can give us sufficient information to find if there are any similarity and common interests between their friends by examining their connections. The pattern of communities between friends of friends also can set rules that help us to spot anomalous users.

**Definition 3 (External Degree).** The *external degree* of $egonet_{i_m}$ to $egonet_{i_n}$ is defined as:

$$d_i(i_m, i_n) = \left| \mathcal{V}^{egonet_{i_m}} \cap \mathcal{V}^{egonet_{i_n}} \right| +$$
$$\left| ij \in \mathcal{E} : i \in \mathcal{V}^{egonet_{i_m}}, j \in \mathcal{V}^{egonet_{i_n}} \right| \quad , i_m, i_n \in \mathcal{G} \tag{3}$$

where $\mathcal{V}^{egonet_{i_m}}$ is the set of nodes of $egonet_{i_m}$ and $\mathcal{V}^{egonet_{i_n}}$ is the set of nodes of $egonet_{i_n}$. The *normalized external degree* is defined as follows:

$$d_i(i_m, i_n)_{\text{norm}} = \frac{d_i(i_m, i_n)}{\min(|i_m|, |i_n|)} \tag{4}$$

**Definition 4 (Community).** The egonets of users $i_m$ and $i_n$ form a *community* if at least the half of the nodes of the smaller egonet connect to the other egonet.

$$\mathcal{C}_i(i_m, i_n) = \begin{cases} 1, & \text{if } d_i(i_m, i_n)_{\text{norm}} \geq \min(|i_m|, |i_n|) / 2 \\ 0, & \text{otherwise} \end{cases} \tag{5}$$

## 2.2 Step 2: Compute Fitting Curve

Local graph metrics related to a single node, its egonet and its super-egonet are used to identify common patterns. We model the relationships between the local metrics using distribution models such as linear and power law. Coefficient of determination ($R^2$) of each model is computed as a goodness of fit measure for the fitting curves.

$R^2 = 1 - SS_{residual} / SS_{total}$ , $SS_{residual} = \sum_{i=1}^{k}(\mathcal{Y}_i - \mathcal{Y}_i^p)^2$, where $\mathcal{Y}_i^p$ is predicted value of $\mathcal{Y}_i$ and $SS_{total} = \sum_{i=1}^{k}(\mathcal{Y}_i - E(\mathcal{Y}_i))^2$, where $E(\cdot)$ gives expected value.

Table 1 includes fitting line equations and $R^2$ for each relationship and dataset; plots are omitted due to page limitations.

**N vs. E (power law) [4]**

Compute a fitting line $E_i \propto N_i{}^a$ , where $1 \le a \le 2$, $E_i$ is the number of edges, $N_i$ is the number of nodes, and $a$ is the power law exponent for user $i$'s egonet.

**E vs. ABC (power law)**

Compute a fitting line $\mathcal{Y} = C\mathcal{X}^\theta$, where $\mathcal{Y}$ is $E$, and $\mathcal{X}$ is $ABC$, and $\theta$ is the power law exponent for user $i$'s egonet.

**N vs. Com (power law)**

Compute a fitting line $\mathcal{Y} = C\mathcal{X}^\theta$, where $\mathcal{Y}$ is $Com$, $\mathcal{X}$ is $N$, and $\theta$ is the power law exponent for user $i$'s super-egonet.

**N vs. E (linear)**

Compute a fitting line $E_i \propto \beta N_i$, where $E_i$ is number of edges, $N_i$ is number of nodes, and $\beta$ is the gradient of the fitting line for user $i$'s egonet.

**E vs. ABC (linear)**

Compute a fitting line $E_i \propto \lambda \sigma_i^{abc}$, where $E_i$ is the number of edges, $\sigma_i^{abc}$ is $\mathcal{ABC}$ and $\lambda$ is the gradient of the fitting line for user $i$'s egonet. Our experiments show there is a relationship between anomaly and the proportion of $E_i$ to $\sigma_i^{abc}$.

## 2.3 Step 3: Compute Outlier Score

For each power law fitting line from step 2, we used the following anomaly score to determine the distance from the fitting line for $ego_i$ ; the calculating follows the OddBall method [4]:

$$aScore(i) = \frac{max(y_i, cx_i^\theta)}{min(y_i, cx_i^\theta)} * \log(|\mathcal{Y}_i - C\mathcal{X}_i^\theta| + 1) \qquad (6)$$

where $\mathcal{Y}_i$ is the y-value, $\mathcal{X}_i$ is x-value of $egonet\ i$, and $\theta$ is a power law exponent. For the power law equation $\mathcal{Y} = C\mathcal{X}^\theta$ this measures "distance to fitting line" by penalizing the number of times that $\mathcal{Y}_i$ deviates from the line.

For each linear fitting line from step 2, we computed $aScore(i)$ in a similar way, but with $\mathcal{Y} = C\mathcal{X} + \theta$ in place of $\mathcal{Y} = C\mathcal{X}^\theta$.

## 2.4 Step 4: Label for Evaluation

Since the existing datasets were not labelled, we used visual inspection to label anomalies. In particular, we visually examined the egonets of each node and decided

whether the node was anomalous our not based on evidence from previous works [1, 3-5]: the majority of users follow the "friends of friends are often friends" pattern and very few users follow either the "cliques or near-cliques" pattern (all the neighbours connected) or the "stars or near-star" pattern (mostly disconnected).

### 2.5 Step 5: Find Threshold

In this step, we compute determine for each metric a threshold value on the outlier score $aScore$ that minimizes the *F-Score*, which is the number of false positives and false negatives in the labelled dataset from step 4. The *F-Score* is calculated as *F-Score =2∗Precision∗Recall / (Precision +Recall)*; its highest value (1) indicates perfect classification of labelled data, whereas its lowest value (0) indicates completely wrong classification of labelled data.

## 3 Experimental Results

Our proposed method is evaluated with three real-life datasets Orkut, Flickr, and Facebook. These datasets were collected by crawling techniques in 2008 [12]. The Orkut dataset has 3M nodes and 23M edges; the Flickr dataset has 1.8M nodes and 22M edges; and the Facebook dataset has 64K nodes and 1.5M edges.

We applied the proposed framework to 20,000 randomly sampled egonets from each dataset. After computing the graph metrics (step 1), fitting curves were computed using regression to determine relationships between metrics (step 2). Outlier scores were then computed for each node (step 3). A labelled subset of 100 nodes (step 4) was then used in threshold finding (step 5) to identify a threshold outlier score that minimizes false negatives and false positives. The resulting F-score was calculated to allow comparison of metrics.

Table 1 compares our observed results for the various graph properties for each of our datasets. We compare five metrics: N vs. E (Linear), E vs. ABC (Linear), E vs. ABC (Power law), N vs. Com (Power law), and N vs. E (Power law); the last being the "OddBall" method of Akoglu et al. [4].

**Table 1.** Comparison of effectiveness graph theoretic properties for anomaly detection in real-life datasets

| Dataset | Method | Fitting curve | $R^2$ | Recall % | Precision % | *F-score* % |
|---------|--------|---------------|-------|----------|-------------|-------------|
| Facebook | E vs. N (Linear) | $V = 0.0638 * E + 29.223$ | 0.80 | 50.51 | 100.00 | 67.11 |
| | E vs. ABC (Linear) | $E = 0.0281 * ABC + 10.553$ | 0.73 | 92.45 | 98.00 | 95.15 |
| | E vs. ABC (Power law) | $E = 0.3839 * ABC^{0.7019}$ | 0.86 | 100.00 | 100.00 | 100.00 |
| | N vs. Com (Power law) | $Com = 0.0369 * N^{2.3508}$ | 0.77 | 52.08 | 100.00 | 68.49 |
| | N vs. E (Power law) [4] | $E = 0.5454 * N^{1.571}$ | 0.95 | 49.49 | 98.00 | 65.77 |
| Flickr | E vs. N (Linear) | $V = 0.009 * E + 187.39$ | 0.65 | 70.00 | 98.00 | 81.67 |
| | E vs. ABC (Linear) | $E = 144.77 * ABC - 8272.1$ | 0.62 | 70.42 | 100.00 | 82.64 |
| | E vs. ABC (Power law) | $E = 0.6151 * ABC^{0.6401}$ | 0.91 | 77.78 | 98.00 | 86.73 |
| | N vs. Com (Power law) | $Com = 0.1248 * N^{2.0304}$ | 0.88 | 57.78 | 52.00 | 54.74 |
| | N vs. E (Power law) [4] | $E = 0.3098 * N^{1.6644}$ | 0.96 | 48.39 | 90.00 | 62.94 |
| Orkut | E vs. N (Linear) | $V = 0.0513 * E + 54.272$ | 0.64 | 77.78 | 75.68 | 76.71 |
| | E vs. ABC (Linear) | $E = 25.025 * ABC + 177.83$ | 0.61 | 94.87 | 100.00 | 97.37 |
| | E vs. ABC (Power law) | $E = 0.544 * ABC^{0.6483}$ | 0.82 | 96.97 | 86.49 | 91.43 |
| | N vs. Com (Power law) | $Com = 0.1045 * N^{2.0689}$ | 0.89 | 75.68 | 75.68 | 75.68 |
| | N vs. E (Power law) [4] | $E = 0.5362 * N^{1.5676}$ | 0.90 | 100.00 | 83.78 | 91.18 |

As we can see from Table 1, our results find that the E vs. ABC (Power law) and E vs. ABC (Linear) methods have the best overall performance across the three datasets. These methods both have higher *F-score* that N vs. E (Power law), which is the Odd-Ball method of Akoglu et al. [4].

## 4     Conclusion

The direct connectivity of online social networks can facilitate illegal activity. In this paper, we have expanded previous research of static analysis of user relationships for detecting anomalous behaviour in online social networks. We have introduced metrics based on a variety of graph properties and presented a framework for detecting nodes with anomalous relationships with other nodes. We applied our approach to datasets from existing online social networks with a manually labelled set of nodes based on existing observations. We identified several metrics—involving the relationship between number of edges and average betweenness centrality of a user's immediate neighbourhood—that perform better than previous. Interesting future work in this area includes the consideration of other datasets and labelling, specifically datasets with pre-established labelling of anomalies, such as email spam or criminal records.

## References

1.  Shrivastava, N., A. Majumder, and R. Rastogi. *Mining (social) network graphs to detect random link attacks*. 2008. IEEE.

2.  Newman, M.E.J., D.J. Watts, and S.H. Strogatz, *Random graph models of social networks.* Proceedings of the National Academy of Sciences of the United States of America, 2002. 99(Suppl 1): p. 2566.

3.  Tong, H. and C.Y. Lin. *Non-negative residual matrix factorization with application to graph anomaly detection*. 2011. SDM.

4.  Akoglu, L., M. McGlohon, and C. Faloutsos, *OddBall: Spotting anomalies in weighted graphs.* Advances in Knowledge Discovery and Data Mining, 2010: p. 410-421.

5.  Sun, J., et al. *Neighborhood formation and anomaly detection in bipartite graphs*. in *Fifth IEEE International Conference on Data Mining*. 2005.

6.  Limsaiprom, P. and P. Tantatsanawong. *Social network anomaly and attack patterns analysis*. in *6th International Conference on Networked Computing (INC)*. 2010.

7.  Heard, N., et al., *Bayesian anomaly detection methods for social networks.* The Annals of Applied Statistics, 2010. 4(2): p. 645-662.

8.  Ball, B., B. Karrer, and M. Newman, *An efficient and principled method for detecting communities in networks.* Arxiv preprint arXiv:1104.3590, 2011.

9.  Yang, Y., Y.C. Guo, and Y.N. Ma, *Characterization of Communities in Online Social Network.* Proceedings of 2010 Cross-Strait Conference on InformationScience and Technology, 2010: p. 600-605799.

10. Clausei, A., C.R. Shalizi, and M.E.J. Newman, *Power-law distributions in empirical data.* Arxiv preprint arxiv:0706.1062, 2007.

11. Brandes, U., *A faster algorithm for betweenness centrality.* Journal of Mathematical Sociology, 2001. 25(2): p. 163-177.

12. Mislove, A., et al. *Measurement and analysis of online social networks*. in *Proceedings of the 7th ACM SIGCOMM conference on Internet measurement*. 2007. San Diego, California, USA: ACM.