

# Uncovering Influence Cookbooks : Reverse Engineering the Topological Impact in Peer Ranking Services

Erwan Le Merrer  
Technicolor, France

Gilles Trédan  
LAAS/CNRS, France

## ABSTRACT

Ensuring the early detection of important social network users is a challenging task. Some peer ranking services are now well established, such as PeerIndex, Klout, or Kred. Their function is to rank users according to their influence. This notion of influence is however abstract, and the algorithms achieving this ranking are opaque. Following the rising demand for a more transparent web, we explore the problem of gaining knowledge by reverse engineering such peer ranking services, with regards to the social network topology they get as an input. Since these services exploit the online activity of users (and therefore their connectivity in social networks), we provide a method allowing a precise evaluation of the extent to which given topological metrics regarding the social network are involved in the assessment of the final user ranking. Our approach is the following : we first model the ranking service as a black-box with which we interact by creating user profiles and by performing operations on them. Through those profiles, we trigger some slight topological modifications. By monitoring the impact of these modifications on the rankings of created users, we infer the weight of each topological metric in the black-box, thus reversing the service influence cookbook.

## INTRODUCTION

The need for an increased transparency in the functioning of web-services has recently arised, motivated by various use cases such as privacy or copyright control. For example, work such as [7] proposes to retrieve which piece of information of a user-profile triggered advertisement to that user. Goal is thus to infer the internals of black-box services provided by companies on the web. Klout or PeerIndex propose to rank users based on their behavior on social networks (using their social connectivity and activity). They nevertheless keep secret the algorithms and parameters used for this ranking<sup>1</sup>. This motivated some users to try reversing their internals [4]. Sometimes information leaks about some ingredients of those hidden recipes ; CEO of PeerIndex admitted to use Pagerank<sup>2</sup> (and thus graph topological-metrics), as a part of the ranking algorithm, to compute user intrinsic influence in a network. Such an understanding of which metrics are involved is also of a particular interest for information sharing and coordination, as it has been shown that some centrality metrics correlate with the actual ability of network actors to coordinate others [5, 3]. This knowledge can then serve to assess if the centrality metric leveraged by the ranking function makes the

ranking service relevant to dispatch roles for given tasks for example [3].

Nevertheless, reverse engineering such black-boxes is a challenging task. Indeed, in this web-service paradigm, the user only has access to the output of the algorithm, and cannot extract any side-information. Moreover, in many cases such as in peer ranking services, the user can only take action on a limited part of the algorithm input. Motivated by this challenge for transparency, we ask the following question : **can a user infer, from the results returned by such peer ranking algorithms, what are the topological metrics in use, and to what extent ?**

We first introduce the ranking service we consider and model our actions, before warming-up on a toy example. We then generalize the example and provide a construction to identify the use of a single arbitrary centrality among a given set of candidates. Then, we assume that the ranking can be produced by a linear combination of multiple centralities, and give a generic reverse engineer approach. We conclude by illustrating such a generic approach on a concrete scenario, before giving perspectives.

## MODEL & WARM UP : REVERSING ONE CENTRALITY

Let us model a social web-service. Each user is represented by a vertex  $v$ , together with a set of (possibly unknown to the user) attributes  $a(v)$ . To interact with the web-service, users have access to a finite set of actions  $A$ . We consider two types of actions : *i*) single actions that only involve a single user (e.g., posting a message on a wall) and might change part of the user profile  $a(v)$ . And *ii*), pair actions that involve a pair of users (e.g., following, declaring or deleting a “friendship” relation). These actions impact the network of relations among users, that we capture as a graph  $G_\infty(V, E)$ .

Among the features of this web-service, a ranking of the users is available. While the internals of the ranking methodology are unknown, each user accesses its output, that is her own ranking at any time. Let  $f$  be the ranking black-box function. More specifically,  $f$  takes as input the graph  $G_\infty$  along with the attributes of its nodes (that is  $\{a(v), \forall v \in V\}$ ) and assigns each node a score  $f(i, a(i))$ ,  $i \in G_\infty$  from which is derived an observable ranking  $r$  of all graph nodes such that :  $\forall i, j \in V(G_\infty)^2$ ,  $i >_r j$  iff  $f(i, a(i)) > f(j, a(j))$ , that is : “node  $i$  is more important (or “influent”) than node  $j$ ”.

The objective of this paper is to gain knowledge on  $f$ , and more specifically to evaluate the impact of each action in  $A$  on users rankings. For a given user, the two main difficulties are that first, she witness only a limited part of the input of  $f$  (typically her own friends in the social graph). Second, the output of  $f$  is sparse, as it only provides nodes with a total

1. Those services may provide a score as an output. Clearly, reversing a ranking function is harder than reversing a score, as you can obtain a ranking from scores, while the opposite is impossible.

2. blog post on <https://www.quora.com/>

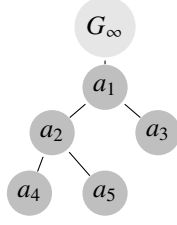


Figure 1: A small query graph  $G_Q$ , solving the single centrality reverse engineering problem for plausible set  $C_{base}$ .

order relation (e.g., user  $x$  is better ranked than her neighbor  $y$ ). In order to try reversing  $f$ , we assume the *querying* user is able to create a set of profiles  $V_a$  in the social service, and have those profiles issue any single action of  $A$ . She is also able to achieve any pair action between two profiles of  $V_a$ , therefore updating the subgraph of  $G_\infty$  induced by nodes of  $V_a$ . Those two operations are conducted through API calls, as it is e.g., observed in practice in Facebook [2].

As a warm up, let us assume that  $f$  leverages exactly one of the following classic centralities  $C_{base} = \{degree, eccentricity, betweenness, Pagerank, closeness\}$  [6]. To determine which one is in use, one user wants to build a small *query graph*  $G_Q$ , attached to  $G_\infty$  (then  $G_\infty \leftarrow G_\infty \cup G_Q$ ), in order to reverse  $f$ . To start our analysis of  $f$  on a clean basis, the user creates nodes  $\in G_Q$  that are strictly identical up to their connectivity (i.e., their attributes in  $a(v)$  regarding single actions such as tweets or posted comments are empty).

**LEMMA 1.** *The query graph  $G_Q$  depicted on Figure 1, of 5 nodes, is sufficient to reverse engineer a function  $f$  that is based on a single centrality  $\in C_{base}$ , relatively to the other centralities in the same set  $C_{base}$ .*

**PROOF.** The proof requires showing that such  $G_Q$  is able to discriminate the centralities considered in the set  $C_{base}$ . Consider graph  $G_\infty \cup G_Q$  on Figure 1.  $G_Q$  nodes are given the following ranking, for centralities in  $C_{base}$ <sup>3</sup>:  $\langle degree, [a_1 =_r a_2 >_r a_3 =_r a_4 =_r a_5] \rangle$ ,  $\langle eccentricity, [a_1 >_r a_2 =_r a_3 >_r a_4 =_r a_5] \rangle$ ,  $\langle betweenness, [a_1 >_r a_2 >_r a_3 =_r a_4 =_r a_5] \rangle$ ,  $\langle Pagerank, [a_2 >_r a_1 >_r a_4 =_r a_5 >_r a_3] \rangle$ ,  $\langle closeness, [a_1 >_r a_2 >_r a_3 >_r a_4 =_r a_5] \rangle$ . All rankings are indeed unique, thus allowing to designate the centrality used, by user observing rankings produced by  $f$  at  $G_Q$  nodes she controls.  $\square$

Note that  $G_Q$  is not the unique graph solving this problem instance.

There are obvious interests in minimizing the size of the constructed query graph : first, constructing a bigger graph requires a longer time, especially if actions on the service platform are rate-limited on operations. Second, the bigger the query, the easier it can be detected by the social service. Note that the graph  $G_Q \setminus n5$ , of size 4 is not a solution, as *degree*

<sup>3</sup>. we conducted numerical simulations using the networkx library : <https://networkx.github.io/>

and *betweenness* produce the same  $[a_1 >_r a_2 >_r a_3 =_r a_4]$  ranking, as for both fringe nodes  $a_3$  and  $a_4$ , *betweenness* is 0, and *degree* is 1.

## GENERAL DISCRETE CENTRALITY DISCRIMINATION

We now generalize the reversing logic used on the previous example to a set  $C$  of arbitrary centralities, possibly in use nowadays. Furthermore, we extend the notion of centrality to the one presented in [1] : a centrality is **any node-level measure**.

We first draw two observations : discrimination is made by the ranking, therefore to distinguish between  $d$  different centralities one requires at least  $d$  different rankings. Thus  $|G_Q| \geq d$ . Second, the discrimination in this set of centralities is made thanks to graphs we call *delta-reversal graphs*.

**DEFINITION 1 (DELTA-REVERSAL GRAPHS).**  $\Delta_{XY}$  is the set of graphs such that  $\forall G \in \Delta_{XY}, \exists i, j \in V(G)$  s.t.  $f_X(G, i) < f_X(G, j) \wedge f_Y(G, i) > f_Y(G, j)$ .

A delta-reversal graph for two centralities  $X$  and  $Y$  is a graph where the ranking  $r$  induced by using the ranking provided by  $f_X$  (i.e., by a function  $f$  only relying on centrality  $X$ ) on the nodes of  $G$  would be different than the ranking induced by  $f_Y$ . Any such graph would thus allow to discriminate between  $X$  and  $Y$  being used as  $f^4$ . The following property is a very handy property for using delta-reversal graphs.

**DEFINITION 2 (CENTRALITY  $k$ -LOCALITY).** Let  $X$  a centrality.  $X$  is said  $k$ -local if  $\forall G_1, G_2$  graphs,  $\forall i \in V(G_1), j \in V(G_2), V^k(i, G_1) = V^k(j, G_2) \Rightarrow f_X(i, G_1) = f_X(j, G_2)$ , where  $V^k(i, G)$  is the graph induced by the  $k$ -hop neighborhood of  $i$  in  $G$ .

The intuition is the following : a  $k$ -local centrality only considers the  $k$ -hop neighborhood of a node when assessing its importance. This can be seen as the “scope” of a centrality : any topological modification beyond this scope leaves the node importance unchanged. This can be exploited to join Delta-reversal graphs into one single query graph while maintaining their discriminating power. Following this intuition, the following definition states an important property of those graphs.

**DEFINITION 3.** Let  $G$  a  $\Delta_{XY}$  graph, and  $dist(i, j)$  the hop-distance between nodes  $i$  and  $j$ . If  $\exists i, j, k \in G$  s.t.  $dist(i, k) > \ell \wedge dist(j, k) > \ell \wedge f_X(G, i) < f_X(G, j) \wedge f_Y(G, i) > f_Y(G, j)$  then  $G$  is  $\ell$ -discriminating.  $k$  is called an anchor.

## Combining Delta-reversal graphs

We now explain how to combine pairwise discriminating graphs into a single query graph.

**LEMMA 2.** Let  $X, Y, Z$  three centralities and let  $k$  their maximum locality. Then  $\forall G_1 \in \Delta_{X,Y}, G_2 \in \Delta_{X,Z}, G_3 \in \Delta_{Y,Z}$ , if all these graphs are  $\ell > k$ -discriminating, then  $G_S = (V(G_1 \cup G_2 \cup G_3) \cup \{a\}, E(G_1 \cup G_2 \cup G_3) \cup \{(a, m_1), (a, m_2), (a, m_3)\}) \in \Delta_{XYZ}$ .

<sup>4</sup>. Examples of discriminating graphs are known in the literature, as they serve as motivation for introducing new centralities : see for instance [8], where a graph is presented that discriminates *random walk betweenness* from classic betweenness centrality.

**Data:**  $G_\infty$ , a target node  $a \in V(G_\infty)$ , the set  $C$  of suspected centralities ( $|C| = d$ ),  $D$  the set of pairwise discriminating graphs for set  $C$

**Result:** The centrality  $X$  in use in  $f$

```

1  $\forall G_{XY} \in D$ , let  $i_{XY}, j_{XY}$  s.t.
    $f_X(i_{XY}) > f_X(j_{XY}) \wedge f_Y(i_{XY}) < f_Y(j_{XY})$ ;
2 //Building and attaching the general query graph to  $G_\infty$ 
3 for  $\forall G \in D$  do
4    $V(G_\infty) \leftarrow V(G_\infty) \cup V(G)$ ;
5    $E(G_\infty) \leftarrow E(G_\infty) \cup E(G) \cup (a, \text{anchor}(G))$ ;
6  $r \leftarrow r(f(G_\infty))$ ;
7 Let  $M$  be a  $d \times d$  matrix initialized to false;
8 //Retrieving the centrality in use in  $f$ 
9 for  $a = 1 \dots d$  do
10  for  $b = a + 1 \dots d$  do
11    $M_{a,b} = i_{X_a X_b} >_r j_{X_a X_b}$ ;
12 Let  $s$  be s.t.  $\forall k = 1 \dots d, M_{s,k} = \text{true}$ ;
13 return  $X_s$ ;

```

**Algorithm 1:** A reverse engineering algorithm, identifying the centrality in use in arbitrary centrality set  $C$ .

PROOF. Since  $G_1$  is discriminating,  $m_1$  exists. Let  $i_1, j_1$  the corresponding anchor nodes. Let  $\sigma_X(G_1), \sigma_Y(G_1)$  the ranks of centralities  $X, Y$ . Assume w.l.o.g. that  $\sigma_X(G_1, i_1) > \sigma_X(G_1, j_1)$  and yet  $\sigma_Y(G_1, i_1) < \sigma_Y(G_1, j_1)$ . Consider  $i_1$ : we have  $d(i_1, m_1) > k$  and thus  $V^k(i_1, G_1) = V^k(i_1, G_S)$ . Thus  $f_X(i_1, G_1) = f_X(i_1, G_S)$ . As the same applies for  $j_1$  we deduce that  $\sigma_X(G_S, i_1) > \sigma_X(G_S, j_1)$  and yet  $\sigma_Y(G_S, i_1) < \sigma_Y(G_S, j_1)$ .

Thus  $G_S \in \Delta_{XY}$ . A similar reasoning holds for  $i_2, j_2$  and  $i_3, j_3$  thus  $G_S \in \Delta_{XZ} \cap \Delta_{YZ} \cap \Delta_{XZ} = \Delta_{XYZ}$ .  $\square$

This lemma is very useful, as it provides us with a way to create discriminating graphs from pair of known ones. The following lemma finally generalizes the construction :

**LEMMA 3.** *Let  $C$  a set of  $d$  centralities and let  $k$  their maximum locality. Let  $D = \{G_{AB} \in \Delta_{AB}, \forall A, B \neq A \in C^2\}$ , a set containing a pairwise discriminating graph for each pair of centrality in  $C$ . If all these graphs are  $\ell > k$ -differentiated, then let  $G_S = (V(\cup_{G \in D} G) \cup \{a\}, E(\cup_{G \in D} G) \cup \{(a, m_{AB}), \forall A, B \neq A \in C\})$ , where  $m_{XY}$  is an anchor of  $G_{XY} \in D$ . Then  $G_S \in \Delta_C$ .*

PROOF. (sketch) : identical to Lemma 2.  $\square$

The  $G_S$  construction therefore allows for any set of centralities, given pairwise discriminating graphs, to construct one general discriminating graph achieving the reverse engineering process. Note that the complexity is quadratic : a graph to compare  $d$  centralities requires  $\Omega(d^2)$  pairwise discriminating graphs.

We are now ready to propose a general method to infer which centrality is in use in  $f$ . It is shown in Algorithm 1.

**THEOREM 1.** *Let  $G_\infty$  a graph, and  $r$  an unknown ranking function relying on centrality  $z$ . If  $z \in C$  then Algorithm 1 returns  $z$ .*

PROOF. First, observe that in Algorithm 1, lines 1 – 5 implement the construction of a combined Delta-reversal graph as defined in Lemma 3. Line 6 collects the resulting ranking. Consider  $M$  at line 12. For  $z$  to be correctly identified, two conditions must hold : *i*) the line  $M_{z,\cdot}$  contains only entries at true, and *ii*) all other lines  $M_{i,\cdot}, i \neq z$  contain at least one false entry.

Consider line  $M_{z,\cdot}$ . By contradiction, assume that one entry, say  $j$  is false. Then necessarily  $i_{z,j} <_r j_{z,j}$  line 11. Since  $r$  is obtained using  $z$ , we deduce  $f_z(i_{z,j}) < f_z(j_{z,j})$ . This contradicts the definition of  $i_{z,j}$  and  $j_{z,j}$  that are chosen line 1 in the sub-graph  $G_{zj}$  such that  $f_z(i_{z,j}) > f_z(j_{z,j})$ . We conclude that  $M_{z,\cdot}$  contains only true entries.

Now, assume there exists another line, say  $i$ , such that  $M_{i,\cdot}$  contains only true entries. Consider column  $z$  : we have  $M_{i,z} = \text{true}$ . As in the previous step, we deduce  $f_z(i_{i,z}) > f_z(j_{i,z})$ ; this again contradicts the definition of  $i_{i,z}$  and  $j_{i,z}$  chosen line 1 in  $G_{iz}$  such that  $f_z(i_{i,z}) < f_z(j_{i,z})$ . Thus every other line has at least a negative entry.

Therefore, we conclude that  $X_s = z$  line 13 : Algorithm 1 has identified  $z$ .  $\square$

The sketch presented in Algorithm 1 can be optimized in many ways. First, one can build the query graph incrementally and only test the relevant centralities : let  $G_{ab}$  be the first added Delta-reversal graph line 4 and 5. It is possible to test directly the value of  $M_{ab}$ . Assume  $M_{a,b} = \text{False}$ , then necessarily centrality  $X_a$  is not used in  $f$ . There is therefore no need to add any other  $G_{ac}, \forall c \in D$  graph.

Second, observe that we focus on pairwise Delta-reversal graphs. Some Delta-reversal graphs allow to differentiate between more than two centralities (for instance, the graph  $G_Q$  Figure 1 that allows to differentiate between 5 centralities at once, while containing only 5 nodes). Using such graphs drastically reduces the size of final the query graph.

## REVERSE ENGINEERING A LINEAR COMBINATION OF CENTRALITIES

In the previous section, we have seen how to identify which centrality is used given a finite set of suspected centralities. We now propose a method for extending to a  $f$  that is a linear combination of suspected centralities, for it allows more complex and subtle ranking functions.

As the space of possible centralities is theoretically infinite, we assume the user takes a bet on a possibly large list of  $d$  centralities in a set  $C$ , that are potentially involved in  $f$ . We will show that our approach also allows to infer the absence of significant impact of a given centrality in set  $C$ , and thus conclude that it is probably not used in  $f$ .

In a nutshell, the query proceeds as follows. The user leverages an arbitrary node  $a$ , already present in  $G_\infty$ . She then creates  $d$  identical nodes (i.e., profiles) and connects them to  $a$ . The ranking of those  $d$  created nodes is thus the same, by construction. She applies to each node a different serie of API calls (i.e., topological operations, attaching them one node for instance). After each serie, ranking of those nodes changes.

**Data:**  $G_\infty$ , a target node  $a \in V(G_\infty)$ , operations  $\{u_1, \dots, u_d\}$

**Result:** An estimate of  $\mathbf{h}$  (i.e., the vector containing the weight of each centrality in  $f$ )

- 1 Let  $\mathbf{k}$  be a vector of size  $d - 1$  initialized to 0;
- 2 **for**  $1 \leq i \leq d$  **do**
- 3     *//attach a query node to target node a, and conduce operations over it*
- 4     Create node  $a_i : V(G_\infty) \leftarrow V(G_\infty) \cup \{a_i\}$ ;
- 5     Add edge  $(a_i, a) : E(G_\infty) \leftarrow E(G_\infty) \cup \{(a_i, a)\}$ ;
- 6     Apply  $u_i(a_i)$ ;
- 7 W.l.o.g.,  $u_d$  is the operation with the highest impact (that is at this step  $\forall j < d, a_j <_r a_d$ ); Reorder otherwise;
- 8 **for**  $i = 1$  to  $d - 1$  **do**
- 9     *//identify operation thresholds*
- 10      $\mathbf{k}_i \leftarrow \max_{x \geq 1} (u_i^x(a_i) <_r u_d(a_d))$ ;
- 11 *//J is the matrix where each element (i, j) is the impact of  $k_i$  applications of  $u_i$  on the  $j^{\text{th}}$  centrality of node  $a_i$ , minus the impact of operation  $u_d$  on  $a_d$ ;*
- 12 Let  $J_{i,j} = c^j(u_i^{k_i}(a_i)) - c^j(u_d(a_d))$ ;
- 13 Set  $J_{d..} = 0^d$ ;
- 14 **return**  $\text{Ker}(J)$  *//find  $\mathbf{h}$  s.t.  $J\mathbf{h} = 0$ , thus is solution to the reverse engineering of  $f$*

**Algorithm 2:** A general reverse engineering algorithm, estimating the linear weight combination of centralities in  $f$ .

Based on those observed changes, she is able to sort the impact of those calls, and thus to describe the impact of one given call by a composition of smaller effect calls. This allows her to retrieve the weights assigned by  $f$  to the  $d$  centralities in set  $C$ , by solving a linear equation system.

Lets consider the following image : imagine you have an old weighing scale (that only answers “left is heavier than right” or vice-versa) and a set of fruits (say berries, oranges, apples and melons) you want to weigh. Since no “absolute” weighing system is available, the solution is to weigh the fruits relatively to each other, for instance by expressing each fruit as a fraction of the heaviest fruit, the melon. One straightforward approach is to directly test how many of each fruit weigh one melon. This is the approach adopted here. However, the problem here is that in general, we are not able to individually weigh each fruit (centrality). Instead, we have a set of  $d$  different fruit salads. This is not a problem if the composition of each salad is known (i.e., the impact of API calls); one has to solve a linear system : there are  $d$  different combinations that are equal, thus providing  $d$  equations.

#### A reverse engineering algorithm

Black-box function  $f$  relies on arbitrary centralities chosen from a set we denote  $C$  of size  $|C| = d$ . Let  $\mathbf{c}_i \in \mathbb{R}^d$  be the  $d$  dimensional column vector representing each of the  $d$  computed centrality values for a node  $i$  in  $G_\infty$ .

We assume that  $f$  is **linear** in all directions (i.e.,  $f$  is a weighted mean of all centralities) :  $\exists \mathbf{h} \in \mathbb{R}^d$  s.t.  $f(i) = \mathbf{c}_i \cdot \mathbf{h}$ . Reverse engineering the topological impact over the final ranking thus boils down to find  $\mathbf{h}$  (and therefore directly obtain  $f$ ).  $\mathbf{h}$  is then the vector of coefficients corresponding to centra-

$$\underbrace{\begin{pmatrix} j_{1,1} & j_{1,j} & \dots & j_{1,n-1} \\ j_{2,1} & j_{2,j} & \dots & j_{2,n-1} \\ \vdots & j_{i,j} = c^j(u_i^{k_i}(a_i)) - c^j(u_d(a_d)) & \ddots & \vdots \\ 0 & 0 & \dots & 0 \end{pmatrix}}_J \underbrace{\begin{pmatrix} h_0 \\ h_1 \\ \vdots \\ h_{d-1} \end{pmatrix}}_{\mathbf{h}} = 0$$

Figure 2: Solution to reverse engineer  $f$ , uncovering  $\mathbf{h}$ .

lities listed in  $C$ . The user performs operations on  $G_\infty$  through API calls, starting from an existing node  $a$ . We assume she is able to find  $d$  different operations denoted  $\{u_1, \dots, u_d\}$ . Consider for instance one operation of that set, noted  $u_1(i)$ , and that simply adds a neighbor to node  $i : G_\infty(V, E), i \xrightarrow{u_1(i)} (V \cup \{a\}, E \cup \{(i, a)\})$ . Such an operation has an impact on  $i$ 's topological role in  $f$ , let  $\mathbf{u} \in \mathbb{R}^d$  be this impact on all centralities in set  $C : \mathbf{c}_i \leftarrow \mathbf{c}_i + \mathbf{u}$ .

Regarding those operations, we assume that : *i*) the user is able to determine the result of each  $u_i$ 's impact on her created node's centrality values (i.e., she computes  $u_i^k(i), \forall i \leq d$ , and where  $k > 0$  is the number of applications of the operation), and *ii*) they are linearly independent : each operation has a unique impact on computed centralities from set  $C$ .

The query proceeds as shown on Algorithm 2, where notations are defined. First, observe that by construction  $\text{rank}(J) \leq d - 1$ . The last operation  $u_d$  is the reference against which we compare other operations. Line 12 records the maximum number of same  $u_i$  operation applications that lead to the same rank (or close) than a single  $u_d$  operation on another node.

Consider a line  $i$  of  $J$  (L.12 Algorithm 2, also represented on Figure 2). Since at the end  $a_i =_r a_d$  (or close), we have  $(c_{a_i} + u_i^{k_i}(a_i))\mathbf{h} = (c_{a_d} + u_d(a_d))\mathbf{h} \pm u_i(a_i)\mathbf{h}$ . Since by construction  $c_{a_i} = c_{a_d}$ , therefore we seek  $\mathbf{h}$  s.t.  $u_i^{k_i}(a_i)\mathbf{h} - u_d(a_d)\mathbf{h} = 0$ . Or matrix notation :  $J\mathbf{h} = 0 : \mathbf{h}$  is in the kernel of  $J$ .

Intuitively, the fact that we get infinitely many solutions ( $\alpha \cdot \mathbf{h}, \forall \alpha \in \mathbb{R}^+$ ) comes from our observation method : we are never able to observe actual scores, but rather rankings. Since multiplying  $\mathbf{h}$  by a constant does not change the final ranking, any vector co-linear to, e.g.,  $\mathbf{h}/\|\mathbf{h}\|$  is a solution.

One important remark is that one cannot formally claim that one centrality metric is not in use in  $C$  with this algorithm. Assume for instance that one centrality, e.g., number of tweets of the considered node, is  $10^3$  times less important than another centrality, e.g., degree. Then we will not be able to witness its effect unless we produce  $10^3$  tweets. And after  $10^2$  tweets, we will only be able to conclude : number of tweets is at least  $10^2$  times less important than degree. One can reasonably assume that such an imbalance in practice means that one service operator will not compute a possibly costly centrality to use it to such a low extent in  $f$ ; this thus makes our algorithm able to discard barely or not used centralities in  $C$ . Finally, we note that with  $\text{cost}(u_i)$  the number of calls issued by operation  $u_i$ , the total number of operations for weighting two centralities in  $C$  is at most  $\text{cost}(u_d) + \sum_{i=1}^{d-1} k_i \cdot \text{cost}(u_i)$ .

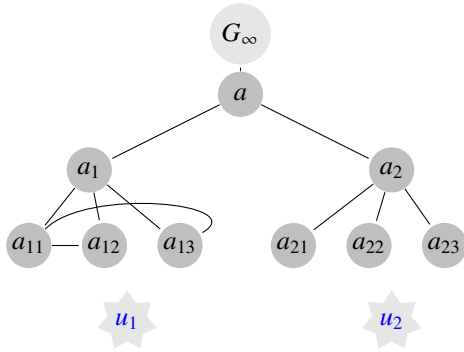


Figure 3: Querying  $G_\infty$  : conducting two sequences of operations  $u_1$  and  $u_2$ , attaching them to  $a$ .

### Exploiting local centralities : an illustration

We demonstrated how to reverse engineer a linear combination of centralities. The difficulty for the user is to compute the impact of  $u$  operations on the suspected centralities. In the easiest case, suspected centralities behave linearly (such as e.g., degree, betweenness), and it is therefore easy to compute the impact of an update. The case of non-linearity can be solved using the locality of centralities : if  $c$  is  $k$ -local, the observation of the  $k$ -hop neighborhood of a node is required to reverse engineer  $f$ . We illustrate this on a simple example.

Let us assume a ranking function  $f$  whose internals use a combination of  $c_1$  : clustering centrality<sup>5</sup> and  $c_2$  : degree (i.e.,  $C = [c_1, c_2]$ ). Without loss of generality, we assume that the coefficient for degree in  $\mathbf{h}$  is  $h_1 = 1$ , so that we seek the corresponding coefficient  $h_2 = h$ . Let us consider the following two operations in  $\{u_1, u_2\}$ . Operation  $u_2$  simply attach a node to one initial query node ( $a_1$  or  $a_2$ ).  $u_1$  starts by attaching  $querySize - 1$  nodes to query node  $a_1$ . At each call, an edge between two randomly selected attached nodes is added, to increase clustering.  $u_1$  and  $u_2$  are represented on Figure 3, for a  $querySize = 4$ . User can compute the value of  $u_1^{k-1}(a_1)$  and  $u_2^{k-1}(a_2)$  at any time, since she controls those nodes.

We simulated the query with a  $G_\infty$  being a 1,000 nodes Barabási-Albert graph with an average degree of 5, estimating  $h$  using  $u_1$  and  $u_2$  operations with Algorithm 1. Figure 4 presents the obtained results : a point  $(x, y)$  means the real value of  $h$  is  $x$  and was estimated by Algorithm 1 as  $y$ . Black dots plot the real coefficient values of  $h$ . Each colored area represents the estimated (reverse engineered) coefficients, while each color represents a query size, i.e., the number of nodes created by the user to reverse  $f$ . The larger the query, the more precise the reverse engineered results. We note that if the real values of coefficients to be estimated are bigger (e.g., 4 or 5 on the  $x$ -axis), estimations show lower precision (larger areas). Despite this remark, estimations appear unbiased.

### DISCUSSION

The will for web-services transparency starts to trigger new research works. XRay [7] for instance proposes a correlation algorithm, that aims at inferring to which data input is associated a personalized output to the user. This Bayesian-based

5. this centrality has no linear behavior, but is 1-local.

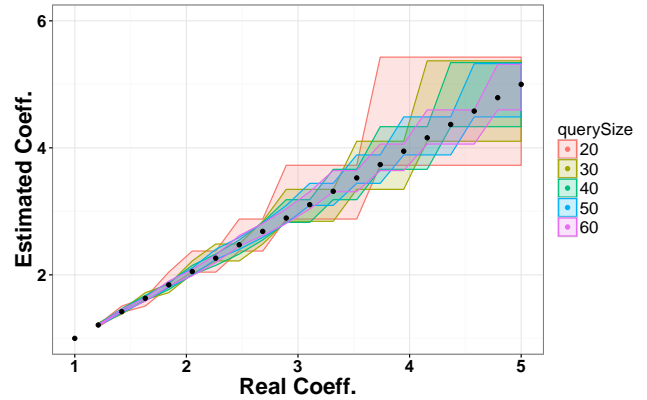


Figure 4: Reversing a  $f$  with unknown coefficients from 1 to 5, with various query sizes (node creations.)

algorithm returns data that are the cause of received ads, while we seek in this paper to retrieve the internals of a black-box ranking function, in order to assess what is the effect of user actions on the output peer ranking. We have presented a general framework. Based on the centralities that might be used by the ranking function, there remain work for a user, for building discriminating graphs, and for finding small topological operations that will make the reverse engineering possible. For a ranking service operator, the countermeasure is the opposite : she must find ranking metrics that are computationally hard to distinguish, typically ones that would ensure the detection of the querying user by the internal security system. We find this to be an interesting challenge for futureworks.

### REFERENCES

1. S. P. Borgatti and M. G. Everett. A graph-theoretic perspective on centrality. *Social networks*, 2006.
2. Q. Cao, M. Sirivianos, X. Yang, and T. Pregueiro. Aiding the detection of fake accounts in large scale social online services. NSDI, 2012.
3. S. Feczak and L. Hossain. Exploring computer supported collaborative coordination through social networks. *The Journal of High Technology Management Research*, 22(2) :121 – 140, 2011.
4. S. Golliver. How i reverse engineered klout score to an  $r^2 = 0.94$ . blog post, 2011.
5. L. Hossain, A. Wu, and K. K. S. Chung. Actor centrality correlates to project based coordination. CSCW, pages 363–372, 2006.
6. D. Koschützki, K. A. Lehmann, L. Peeters, S. Richter, D. Tenfelde-Podehl, and O. Zlotowski. *Network Analysis : Methodological Foundations*, chapter Centrality Indices. Springer, 2005.
7. M. Lécuyer, G. Ducoffe, F. Lan, A. Papancea, T. Petsios, R. Spahn, A. Chaintreau, and R. Geambasu. Xray : Enhancing the web’s transparency with differential correlation. USENIX Security Symposium, 2014.
8. Newman. A measure of betweenness centrality based on random walks. *Social Networks*, 27(1) :39–54, 2005.