# A framework for validating the merit of properties that predict the influence of a twitter user

Stefan Räbiger [a,1], Myra Spiliopoulou [b,*]

[a] Faculty of Engineering and Natural Sciences, Sabanci University, Orhanli, Tuzla, 34956 Istanbul, Turkey
[b] Faculty of Computer Science, Otto-von-Guericke University Magdeburg, Universitätsplatz 2, D-39106 Magdeburg, Germany

## ABSTRACT

What characterizes an influential user? While there is much research on finding the concrete influential members of a social network, there are less findings about the properties distinguishing between an influential and a non-influential user. A major challenge is the absence of a ground truth, on which supervised learning can be performed. In this study, we propose a complete framework for supervised separation between influential and non-influential users in a social network. The first component of our framework, the *InfluenceLearner*, extracts a Relation Graph and an Interaction Graph from a social network, computes network properties from them and then uses them for supervised learning. The second component of our framework, the *SNAnnotator*, serves for the establishment of a ground truth through manual annotation of tweets and users: it contains a crawling mechanism that produces a batch of tweets to be annotated offline, as well as an interactive interface that the annotators can use to acquire additional information about the users and the tweets. On this basis, we have created a ground truth dataset of Twitter users, upon which we study which properties characterize the influential ones. Our findings show that there are predictive properties associated with the activity level of users and their involvement in communities, but also that writing influential tweets is not a prerequisite for being an influential user.

© 2014 Elsevier Ltd. All rights reserved.

## 1. Introduction

The propagation of influence in online social networks has been subject of extensive research, ever since the seminal works of Domingos and Richardson (2001) and Kempe, Kleinberg, and Tardos (2003) on influence propagation in social graphs. Whilst there is a substantial amount of work in identifying influential social graph participants (also known as influentials), there are less findings on identifying the properties which distinguish between influential and non-influential nodes. In this study, we investigate to what extend supervised analysis of a social graph can reveal distinctive properties of influential users. We propose a framework that extracts user attributes that have the potential of predicting a user's influence power, and we use this framework to separate between influential and non-influential Twitter users.

Modeling the spread of influence in social networks is an intensively studied task. Contributions include diffusion models that describe influence propagation, and theoretical findings on how well such a diffusion model can describe reality (Kempe et al., 2003). A major application area for such models is viral marketing, because, as pointed out by Barbieri, Bonchi, and Manco (2013): "…individuals tend to adopt the behavior of their social peers". They continue with the important statement that "cascades happen first locally, within close-knit communities, and become global "viral" phenomena only when they are able cross the boundaries of these densely connected clusters of people." (Barbieri et al., 2013). These studies focus on modeling the spread of influence and on finding the persons that have the most influence. But what are the *properties* characterizing these persons? Are there attributes on the activities and writings of a user that indicate her influence?

In this study, we express the problem of identifying influential users as a classification task, and aim to identify the characteristics of such users. Labeled datasets for this task are rare (Bigonha, Cardoso, Moro, Gonçalves, & Almeida, 2012). Moreover, the authors of Bigonha et al. (2012) are not allowed to give access to the tweet contents due to the terms of service of Twitter. To verify the assumption that tweet content can indicate whether a tweet author is influential, we propose as part of our framework a workflow for the offline creation of a labeled set of tweets and of users

---

* Corresponding author.
   *E-mail addresses:* stefan@sabanciuniv.edu (S. Räbiger), myra@iti.cs.uni-magdeburg.de (M. Spiliopoulou).
   [1] Work done while being a master student at the Faculty of Computer Science, Otto-von-Guericke Univ. Magdeburg.

who wrote these tweets. We use this workflow to create a labeled dataset that we use for our experiments.

The contribution of our work is then twofold: we propose a supervised learning method and a set of properties for distinguishing between influential and non-influential social graph members, and we also propose a workflow for acquiring labeled data for supervised learning. We study Twitter in our work, because Twitter is a representative of the *who listens to whom* attitude suggested in Bakshy, Hofman, Mason, and Watts (2011). However, the workflow we propose allows for building datasets on any social platform to learn a dedicated model on it. As a further by-result of our approach, we make the dataset of our first run of our framework available to other scholars (see Section 4 on data access).

The remainder of this paper is structured as follows: in Section 2 we discuss related work for detecting influentials in Twitter. In Section 3, we introduce our framework and present its learning component that extracts attributes from the social graph to characterize the users of postings, as well as the annotation component for acquiring a labeled dataset. We report on our experiments in Section 4. The last section concludes our study.

## 2. Related work

There exist various heuristics – based on mentions, replies, followers and followees – that rank users according to their influence (Anger & Kittl, 2011; King et al., 2013; Razis & Anagnostopoulos, 2014; Sun & Ng, 2013). Other heuristics focus on aspects like tweet quality (Kong & Feng, 2011) or utilize alpha centrality which is related to Eigenvector centrality (Overbey, Paribello, & Jackson, 2013). Zhao et al. devise a new measure for influence based on sentiment (Zhao et al., 2014). Sun et al. pursue a more sophisticated approach by building a user and tweet graph to identify influential users (Sun & Ng, 2013). King et al. (2013) devise the t-index that denotes the number of times a user's unique tweet has been retweeted to compare the influences users exert on the same topic. Razis and Anagnostopoulos (2014) combine the ratio of a user's followers and followees and the ratio of tweets written in a certain period of time into an influence metric. Similarly, Bigonha et al. combine users' sentiment, tweet quality and centrality to obtain an aggregated influence score (Bigonha et al., 2012). In terms of supervised learning, Chai et al. follow a similar approach, but combine attributes related to four categories – activity, centrality, quality and reputation (Chai, Xu, Zuo, & Wen, 2013). Liu et al. extract several attributes known from literature to train an SVM (Liu, Li, Xu, & Yang, 2014) and Xiao et al. use attributes related to three different categories in order to find influential users (Xiao, Zhang, Zeng, & Wu, 2013). The problem with these approaches is that there is no ground truth on what people consider as influential user to evaluate the approaches on. We address this aspect in our work by building a ground truth on people's perception of influence.

There are also commercial services (including Klout,[2] PeerIndex,[3] Kred[4]) that assign influence scores to users. However, each such service uses its own, internal/proprietary definition of the term influence. Campo-Ávila, Moreno-Vergara, and Trella-López (2013) attempt to reverse engineer two of these algorithms (Klout and PeerIndex) and to identify the factors used in these internal definitions of influence. Our intention in this work is not to provide yet another definition of influence (which might be subject of some controversy), but to identify factors that are associated with influence, *when human annotators decide who is influential and who is not*, keeping in mind that humans, in contrast to services, do not have a rigid definition of whom they consider influential.

In a different thread of research, Barbieri et al. study the spread of information in a social network, and point out that cascades are local phenomena (Barbieri et al., 2013) that manifest themselves inside close-knit communities; only some of them cross the community borders through nodes that are part of both communities. Wang et al. also assume that influentials need to occur in every community to propagate information across the network (Wang, Cong, Song, & Xie, 2010). Inspired by these findings on the role of communities for information propagation, we also take the community structure of the graph into account. However, in contrast to Wang et al. (2010) and Barbieri et al. (2013), our objective is to find the characteristics of influential users and not to point to those users who are influential.

Summarizing, our work differs from other literature on influence in following aspects. Differently from Bigonha et al. (2012), Anger and Kittl (2011), King et al. (2013), Sun and Ng (2013), Razis and Anagnostopoulos (2014), Kong and Feng (2011), Overbey et al. (2013) and Zhao et al. (2014) and similarly to Quercia, Stillwell, Michal Kosinskil, and Crowcroft (2011), we do not attempt to find influential users but rather identify the characteristics that separate between influential and non-influential users. To this purpose, we derive properties that reflect social activity, and use them in supervised learning. The learner and the set of properties constitute our first contribution. This set of properties is larger and more elaborate than in Bigonha et al. (2012), Chai et al. (2013), Liu et al. (2014) and Xiao et al. (2013) and, moreover, it is accompanied by an elaborate approach on assessing the ground truth. Indeed, unlike Klout, PeerIndex, Kred and Bigonha et al. (2012) and Cha, Haddadi, Benevenuto, and Gummadi (2010), we do not provide yet another definition of influence, nor try to reengineer existing definitions, but we rather cover the non-crisp, subjective perception of influence that people have. To this purpose, our approach encompasses a mechanism for the creation of a ground-truth dataset (a seed) of influential and non-influential users through human annotators. This mechanism is our second contribution.

## 3. Framework

Our framework for identifying characteristics of influentials has two components: the *SNAnnotator* and the *InfluenceLearner*. The former collects data from Twitter to establish a ground truth, which serves as input for the *InfluenceLearner*. The latter is responsible for turning a dataset into graphs, extracting attributes and learning a meaningful model.

### 3.1. SNAnnotator

Our *SNAnnotator* describes the process of collecting a dataset regarding a specific topic in batch mode, manually labeling and preparing it for attribute extraction with our *InfluenceLearner*.

### 3.1.1. Offline dataset crawl

We collect tweets during a certain period of time and retrieve their authors thereafter. *SNAnnotator* can also operate on multiple topics simultaneously, because it uses hashtags to identify the tweets corresponding to a topic. However, as Cha et al. point out (Cha et al., 2010), a user's influence may vary over topics and change with time. Therefore, we concentrate on learning influence towards single topics. We collect tweets containing the set of predefined hashtags using the Twitter Streaming API.[5] This means only the latest tweets of users related to the topic are collected. Once this process is completed, the metadata of the respective users are

---

[2] www.klout.com (10-30-2014).
[3] www.peerindex.com (10-30-2014).
[4] www.kred.com (10-30-2014).

[5] https://dev.twitter.com/docs/streaming-apis (10-30-2014).