

Disentangling Sources of Influence in Online Social Networks

MATIJA PIŠKOREC^{1,2}, TOMISLAV ŠMUC¹, AND MILE ŠIKIĆ^{2,3}

¹Rudjer Boskovic Institute (RBI), 10000 Zagreb, Croatia

²Faculty of Electrical Engineering and Computing, University of Zagreb, 10000 Zagreb, Croatia

³Genome Institute of Singapore, Agency for Science, Technology and Research (A*STAR), Singapore 138632

Corresponding authors: Matija Piškorec (matija.piskorec@irb.hr) and Mile Šikić (mile.sikic@fer.hr)

This work was supported in part by the Centre of Excellence Project “DATACROSS”, and in part by the Croatian Government and the European Union through the European Regional Development Fund - the Competitiveness and Cohesion Operational Programme under Grant KK.01.1.1.01.0009.

ABSTRACT Information propagation in online social networks is facilitated by two types of influence - *endogenous* (peer) influence that acts between users of the social network and *exogenous* (external) that corresponds to various external mediators such as online news media. However, inference of these influences from data remains a challenge, especially when data on the activation of users is scarce. In this paper we propose a methodology that yields estimates of both endogenous and exogenous influence using only a social network structure and a single activation cascade. Our method exploits the statistical differences between the two types of influence - endogenous is dependent on the social network structure and current state of each user while exogenous is independent of these. We evaluate our methodology on simulated activation cascades as well as on cascades obtained from several large Facebook political survey applications. We show that our methodology is able to provide estimates of endogenous and exogenous influence in online social networks, characterize activation of each individual user as being endogenously or exogenously driven, and identify most influential groups of users.

INDEX TERMS Data collection, information diffusion, maximum likelihood estimation, social network services, online social networks, statistical learning.

I. INTRODUCTION

Popularity of online social networks allows us to investigate dynamics of social interactions on a scale that was previously unattainable [1]–[8], while at the same time raising ethical concerns not previously encountered [9], [10]. One particular type of social interaction is an *information cascade* - a spread of information between individuals in a social network [11], [12]. Information cascades are instrumental in investigating *social influence*, which can be defined as the degree to which the behavior of individuals changes the behavior of their peers [13]. Although mathematical modeling of social influence and information cascades is an active field of research in sociology for decades [11], [12], it only recently became technologically feasible to apply it to wide range of domains such as viral marketing [14],

information diffusion [15], behavior adoption [16] and epidemic spreading [17].

Presence of *exogenous* factors is particularly problematic in estimation of social influence as it confounds with the *endogenous* factors, and can be hard to differentiate using observational data alone [18]. Still, it is instrumental for understanding the information spreading as information can propagate through multiple channels simultaneously, many of which are exogenous to the online social network itself - news media websites, direct communication via email and instant messengers, and even offline word-of-mouth transmission. In addition, external events such as political unrest [1], [19] and natural disasters [20] are often strong mediators of information cascades. These exogenous influences are usually not directly observable in the online social network itself, although they can be inferred from the available data. Understanding how endogenous and exogenous forces influence the information diffusion in online social networks could help us estimate to what extent are these vulnerable to manipulation

The associate editor coordinating the review of this manuscript and approving it for publication was Xiao Liu.

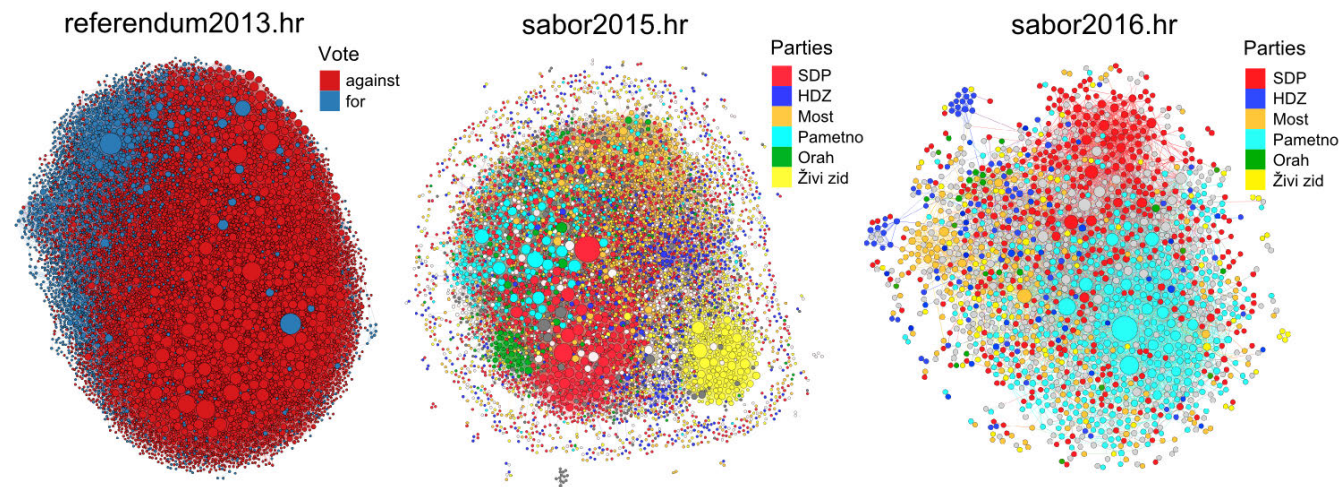


FIGURE 1. Collected Facebook friendship networks. Visualization of Facebook friendship networks of users who registered on three of our Facebook online survey applications: referendum2013.hr (11538 registered users), sabor2015.hr (6909 registered users) and sabor2016.hr (3818 registered users). Network nodes are colored according to the user's votes, and node sizes correspond to the number of their Facebook friends that also registered on the survey application. Clustering of users into communities based on votes shows a homophily effect - users are more likely to associate with other users that share their political preferences. This suggests a potential for endogenous influence.

by various interest groups such as organized individuals, news media and government agencies [21].

In this paper we present a new methodology for estimation of endogenous and exogenous influence in online social networks. Our current model is conceptually similar to the unified model of social influence [22] which was shown to be generalization of many popular influence models, including complex contagion model [4], independent cascade model [23] and generalized threshold model [23]. In our previous work [24] we proposed a simpler method for inference of endogenous and exogenous influence that exploits statistical differences between the way the two types of influence act on users. The underlying assumption is that the endogenous influence is dependent on the current state of the social network and which users are already active or not, while the exogenous influence is independent on these. By incorporating these assumptions in a statistical model we can infer magnitude of endogenous and exogenous influence from empirical data.

Here, we develop a likelihood-based approach which is expressive enough to accommodate many different microscopic models of influence, and propose a maximum likelihood inference method to estimate the parameters. The inference problem is the following - given a single activation cascade and a friendship network between users, and assuming a particular form of endogenous influence, infer parameters of endogenous and exogenous influence and estimate magnitudes of these influences in time and on a global and user level.

We evaluate our methodology on activation cascades collected via three online survey applications related to three distinct political events in Croatia (Fig 1 and Fig 2, and Table 1). First survey, which is related to the referendum on the definition of marriage in 2013, we already used in our previous work [24]. Other two surveys are related to Croatian

TABLE 1. Collected online survey data. Collected online survey data include demographic information, friendships between users, and referral links through which users visited our applications. Time period refers to the period when surveys were active. Depending on whether these referral links originated within Facebook or some external website they could be used as indicators of endogenous and exogenous influence respectively.

Dataset	Time period	Users	Collected data
referendum2013.hr	25.11. - 1.12.2013.	10175	friendships, demographics
sabor2015.hr	2.11. - 8.11.2015.	6909	friendships, referral links
sabor2016.hr	5.9. - 11.9.2016.	3818	friendships, referral links

parliamentary elections in 2015 and 2016 and we collected them exclusively for this research. In all of our surveys the activation cascades are a series of user registrations through time. Surveys were active one week prior to actual elections and through them users were able to express their vote on the upcoming elections, see summary statistics for all users as well as for their online peers, and share the link to the survey through Facebook. Besides votes, we also collected Facebook friendship connections between all users that participated in our survey. In 2013 survey we also collected demographic data and in other two we obtained referral links through which users visited our survey website. These referral links originate either from Facebook, which indicates endogenous influence, or from some external website, which indicates exogenous influence. This classification of referral links served as a proxy for ground truth influence and allowed us to evaluate our inference method. During data collection we followed Facebook's Platform Policy which provides guidelines and regulations for the usage of Facebook Graph API in third-party Facebook applications.

The main contributions of this paper are the following:

- 1) We collected data on social engagement of over 20 thousand Facebook users that participated on three

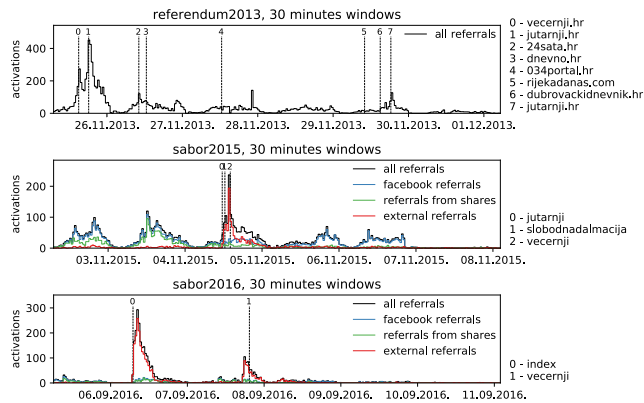


FIGURE 2. Collected registration time series. Collected registration time series of users that registered on our three online survey application. Time series are annotated with times of major news events which reported on our online survey application, and which are used as a proxy for exogenous influence. Time series for sabor2015.hr and sabor2016.hr datasets are additionally separated based on the type of the referral links.

distinct online political surveys. Datasets where users have to provide an informed consent to collect their data are usually much smaller, and so researchers have to rely on simulated datasets in order to validate their models.

- 2) We estimate magnitude of endogenous and exogenous influence in social networks by using only a single activation cascade of users and their friendship network. Most previous research relies on the availability of multiple information cascades and rarely tackles exogenous influence directly by either leaving it as an option [22], devising experiments where it is negligible [25] or simply treat it as a nuisance [26].
- 3) We show how can our methodology be used to estimate collective influence of various groups of users and characterize to what extent was their activation endogenously or exogenously driven. These estimates agree with both the simulated activation cascades and three realistic use cases where user's referral links served as a proxy for the ground truth labels on whether users were endogenously or exogenously activated.

II. RELATED WORK

The most commonly used information diffusion models were inspired by epidemiology which model how a disease spreads in a population [27]–[29]. However, their utility is sometimes hindered by their use of latent states which are unobservable in data. For this it is more appropriate to use Independent Cascade (IC) model [30] and Linear Threshold (LT) model [11], [23] which feature two observable states - *active* and *inactive* that denote whether an user was already exposed to the piece of information or not. These are popular for their simplicity that facilitates theoretical analysis [31], statistical inference from data [26], and can also be used as building blocks for more complex applications such as influence maximization [32], [33]. In this work we use two variants of the IC model (Eq 1 and 2).

However, there are several crucial differences between epidemic spreading and information diffusion [34]. Epidemic spreading is better modeled with simple contagion model where endogenous factors play a dominant role, and the activation probabilities are independent of the neighborhood structure and the state of activated users in it. On the other hand, information diffusion is better modeled with complex contagion due to the common presence of exogenous factors [25] and more complex forms of endogenous influence which include various social reinforcement mechanisms such as reciprocity [35], social feedback [36] and homophily [37]. These additional factors are often neglected in modeling, which is reasonable if there is enough evidence that some of them, for example exogenous influence, is negligible [25]. When this is not possible the exogenous influence has to be explicitly accounted for [38]–[40]. In our work we rely on an explicit modeling of exogenous influence through a likelihood-based approach.

Many likelihood-based approaches for modeling influence exist in literature, including peer and authority model [41] which, however, requires explicit modeling of *authorities* responsible for exogenous influence, while in our case this is not necessary. Many of the other approaches rely on the availability of multiple activation cascades [39], while we use only one. Also, we use the social network structure, based on final state of activation cascade, directly in our inference rather than using it implicitly [8] or relying on a network statistic such as degree distribution [42].

III. METHODS

Crucial components of our methodology are explicit microscopic models of endogenous and exogenous influence with which we expand the Independent Cascade (IC) model. We then use these models in a log-likelihood function which gives us probability of observing particular activation cascade as a function of the model's parameters. Formulating our inference problem in a probabilistic way allows us to optimize for the maximum likelihood parameters and to estimate the magnitude of endogenous and exogenous influence. We apply our methodology on several simulated and empirical activation cascades in order to characterize the activation of users as being more endogenously or exogenously driven. The simulated case is easier because we know both the functional form and the parameters of the model that generated simulated information cascade, which allows us to perform evaluation in a straightforward manner. For the empirical cases we use three Facebook datasets obtained from an online political survey applications. In the end we estimate collective influence of three groups of users - those who registered by following link from within Facebook, those that registered by following link from an external website, and those that followed a link from a Facebook advertisement.

A. MODELS OF ENDOGENOUS AND EXOGENOUS INFLUENCE

We assume that an activation of an user in an online social network is mediated by two influences (Fig 3): (i) endoge-

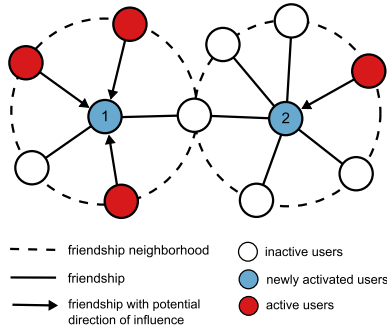


FIGURE 3. Endogenous and exogenous influence. Our assumption is that information propagation in an online social network is mediated by two types of influence - endogenous (peer) which acts between the users of the social network and exogenous influence which is external to it. The estimated endogenous influence on the newly activated user $i = 1$ should be higher because more of his peers are already active, as compared to user $i = 2$.

nous influence p_{peer} which depends on the network structure and users that are already active or not, and (ii) exogenous influence p_{ext} which is modeled as a time dependent random variable and is constant across all users. An additional assumption is that parameters of endogenous influence are constant throughout the period of observation, while parameters of exogenous influence may change in time. Both sets of parameters are equal for all users. This allows us to use a very simple model for the exogenous influence - a single probability of activation $p_{ext}^{(i)}(t)$ which is equal for all inactive users i at each specific time step, although it can change in time. Instead of parameterizing $p_{ext}^{(i)}(t)$ with a suitable closed form, we chose to evaluate it at each time step independently [39]. The benefit of such *nonparametric* estimate of exogenous influence is that it minimizes assumptions on the influence's functional form. For example, we do not have to explicitly incorporate decay of exogenous influence into our model as we will do with the endogenous influence later on (Eq 2), but our model is still able to infer this decay if it is supported by available data. An example of this is visible in the experiment in Fig 6.

For the endogenous influence we choose two commonly used Independent Cascade (IC) models: (i) Susceptible-infected (SI) model $p_{SI}^{(i)}(t)$ and (ii) Exponential decay (EXP) model $p_{EXP}^{(i)}(t)$. IC models are an example of *simple contagion* - activation of users happens due to a direct influence of one of their peers, independently of the rest of the system, including the neighborhood structure and which other users are active or not. EXP model has an added condition that peers that activated recently carry more influence than the ones that activated farther away in time, which is commonly incorporated in endogenous influence models [43], [44].

Probability of endogenous activation for user i at time interval $[t - \Delta t, t]$ under the SI model is defined as follows:

$$p_{SI}^{(i)}(t) = 1 - \prod_{j \in N^{(i)} \text{ active at } t} (1 - p_0) = 1 - (1 - p_0)^{a_i(t)} \quad (1)$$

where $N^{(i)}$ is a set of peers of user i , $a_i(t)$ designates how many of them are active at time t , and p_0 is a probability

of user i 's being activated by each of its peers. Assumption of the SI model is that probability of activating one's peers does not change in time, so once user is activated, every subsequent step he has the same probability p_0 of activating any of his peers. This assumption is more appropriate in epidemiological setting, from where SI model originated, than in information propagation setting where we would expect the influence to decay in time. This could be achieved by adding a parameter for influence decay, which leads us to the EXP model:

$$p_{EXP}^{(i)}(t) = 1 - \prod_{j \in N^{(i)} \text{ active at } t} (1 - p_0 e^{-\lambda(t-t_j)}) \quad (2)$$

where t_j is the time of activation of user j . p_0 and λ are parameters of endogenous influence which define the shape of exponential decay of influence, with p_0 being the probability of user j activating user i at time $t = t_j$ and λ being the half-decay of influence. Both SI and EXP models feature independent cascades - each individual user can independently activate any of his peers. However, in social contagion it is more realistic to add a requirement of multiple interactions for the activation. This effectively models social reinforcement mechanism which is a known driving force for product adoption [25]. One of the simplest examples of such *complex contagion* models is the *threshold model* where the probability of endogenous activation is related to the number of already active peers $N^{(i)}$ of user i . We define one such threshold model in the Eq S11 of the Supporting information and show that it can also be effectively incorporated into our inference methodology.

We now define a likelihood function \mathcal{L} which gives us probability of observing data D (network and activation times) at a particular time t given some functional forms for endogenous and exogenous influence p_{peer} and p_{ext} . Due to typically small probabilities involved in these processes we actually use log-likelihood for maximum likelihood estimation of parameters, where product of probabilities is replaced with the sum of log-probabilities:

$$\begin{aligned} \log \mathcal{L}(D; p_{peer}, p_{ext}, t) &= \sum_{i \in \text{activated at } [t-\Delta t, t]} \log(1 - (1 - p_{peer}^{(i)}(t))(1 - p_{ext}(t))) \\ &+ c(t) \sum_{i \in \text{inactive at } t} \log((1 - p_{peer}^{(i)}(t))(1 - p_{ext}(t))) \end{aligned} \quad (3)$$

First term on the right-hand side quantifies the agreement for the users that *did* activate in a given time period $[t - \Delta t, t]$, as this had to be due to either endogenous or exogenous influence. Second term quantifies the agreement for the users that *did not* activate up to time t , neither through endogenous nor through exogenous influence. The time enters our inference *only* through the activation time of users and is used in two ways - i) to determine which users were active or inactive in time window $[t - \Delta t, t]$ (Eq 3), and ii) to calculate endogenous influence decay in EXP model (Eq 2). However, in principle it is possible to use a *temporal network* where

friendship connections between users change in time. This can be encoded into the expression for endogenous influence p_{peer} , for example in equations for the SI and EXP models (Eq 1 and 2) by making $N^{(i)}$ - a set of peers of user i , a time-changing quantity. We can remove explicit dependence on time t from the Eq 3 by evaluating \mathcal{L} nonparametrically - at each time increment Δt .

One issue still needs to be addressed - on which users does the exogenous influence actually acts? We know that our friendship network does not contain *all* possible users, and so the true number of yet inactive users is probably much larger than what we actually observe. This *observer bias* could lead to the overestimation of the exogenous influence as we approach the end of the activation cascade and the number of eventually observed inactive users decreases towards zero, while the true number of inactive users which could possibly activate (but do not during our observation period) stays large. We correct for this by artificially increasing the part of our log-likelihood which is responsible for inactive users by factor $c(t) = 1 + \alpha(N_{all}/N_{inactive}(t))$, where N_{all} is the number of all users in the social network, and $N_{inactive}(t)$ the number of all yet users inactive users at time t (more details in Section S6 of Supporting information).

B. ALTERNATING METHOD FOR INFERENCE

Our two main assumptions during statistical inference are: (i) both endogenous and exogenous influence are equal for all users at any given time, and (ii) endogenous influence does not vary in time while exogenous influence does. This leads us to the inference algorithm where we seek a single set of parameters for the endogenous influence p_{peer} and a set of parameters for the exogenous influence $\{p_{ext}\}_t$ for each time step t . This would make the dimensionality of our log-likelihood proportional to the number of time steps we use for inference, which would be hard to optimize numerically. Instead, we use an *alternating* method [39] where we alternately fix either p_{peer} or $\{p_{ext}\}_t$ and optimize for the other.

Algorithm 1 gives the pseudocode of the alternating procedure for inference of endogenous p_{peer} and exogenous $\{p_{ext}\}_t$ influence that we use in our experiments. In the first part of the algorithm (steps 2-4) we optimize p_{peer} and p_{ext} for every time window separately, which then serve as initial values for the alternating procedure. Optimization procedure is designated with a generic *MAP* (Maximum A Posteriori) procedure which takes as arguments the parameters which are held fixed and outputs values of the remaining parameters so that the log-likelihood (Eq 3) is maximized. The actual *MAP* optimization is performed with a truncated Newton algorithm that is Hessian-free and uses conjugate gradients to iteratively compute parameter updates [45], although in principle any suitable optimization algorithm could be used. Second part of the algorithm is the actual alternating procedure (steps 5-11) where we first optimize for a single set of endogenous parameters p_{peer} , conditioning on the exogenous parameters $\{p_{ext}\}_t$ we obtained for each time window (step 6). We then optimize exogenous parameters for each window

separately $\{p_{ext}\}_t$, conditioning on a single set of endogenous parameters p_{peer} we obtained in the previous step (step 7). We then alternate between the step 6 and 7 until values for p_{peer} and $\{p_{ext}\}_t$ converge. The difference between the values for the current and previous iteration are calculated in steps 8 and 9 and the convergence itself is checked in step 5.

C. INFERENCE OF ACTIVATION TYPES

Because our model gives us probabilities for endogenous and exogenous activation for each user individually, we can use this information to estimate activation type for each of the users. For this we define a single measure of *exogenous responsibility* $R^{(i)}$ which quantifies to what degree is an activation of user i due to the exogenous (external) influence:

$$R^{(i)}(t) = \frac{p_{ext}(t)}{p_{ext}(t) + p_{peer}^{(i)}(t)} \quad (4)$$

where t is the time of activation of user i . Values close to zero indicate dominating endogenous influence, and values close to one indicate dominating exogenous influence. An extreme value of zero is achieved for users who activated during time when there was no exogenous influence acting in the network. An extreme value of one is achieved for users who, at the time of their activation, did not have any active peers. Note that it is not possible for both $p_{ext}(t)$ and $p_{peer}^{(i)}(t)$ to be 0, and consequently that the value of responsibility is undefined, because that would mean the activation of this user is evaluated as *impossible* by our model in Eq 3. In principle, we could also use pure activation probabilities $p_{peer}^{(i)}$ or p_{ext} as measures of influence, but experiments on simulated data showed that exogenous responsibility is the most sensible (more details in Supporting information).

D. INDIVIDUAL AND COLLECTIVE INFLUENCE OF USERS

Our assumption is that each user is, to some extent, responsible for endogenous activation of all of his peers that activated after him. This influence extends beyond user's immediate peers. However, as we do not have a deterministic activation path (we do not know who shared information with whom) it is not straightforward to transitively incorporate influence from far away users as it is usually done [46]. This is why we express the influence $I^{(i)}$ of user i (Eq 5) as the extent to which user i is responsible for activation of his peers j :

$$I^{(i)} = \sum_{j \in N^{(i)}} \frac{I^{(i \rightarrow j)}}{\sum_{m \in N^{(i)}} I^{(m \rightarrow j)}} p_{peer}^{(j)}(t_j) \quad (5)$$

where $I^{(i \rightarrow j)}$ is the fraction of the endogenous influence that user i can claim for user j . In our case we define it as $I^{(i \rightarrow j)} = 1$ if i and j are peers, and 0 otherwise. This means that all user's are credited equally for the activation of their peers, regardless of how far away in time they themselves activated. For an alternative formulation which involves time see Eq S8 in the Supporting information. As shown on Fig 4, each user can claim part of the peer activation probability $p_{peer}^{(j)}(t_j)$ for each of his peers j that activated after him $t_i < t_j$. As we do not have a deterministic activation path, this is really

Algorithm 1 Alternating Method for Joint Inference of Influence

```

1: procedure AlternatingInference( $T, \epsilon, p_{peer}(t), p_{ext}(t)$ )
2:   for  $t \in \{1, \dots, T\}$  do
3:      $\{p_{peer}\}_t, \{p_{ext}\}_t \leftarrow \text{MAP}(p_{peer}(t), p_{ext}(t))$  ▷ Optimize for every time window.
4:   end for
5:   while  $\Delta_{peer}^{(i-1)} \geq \epsilon$  &  $\Delta_{ext}^{(i-1)} \geq \epsilon$  do ▷ Until  $p_{peer}$  and  $\{p_{ext}\}_t$  converge.
6:      $p_{peer}^{(i)} \leftarrow \text{MAP}(\{p_{ext}^{(i-1)}\}_t)$  ▷ Fix  $\{p_{ext}^{(i-1)}\}_t$  and optimize for single  $p_{peer}^{(i)}$ .
7:      $\{p_{ext}\}^{(i)} \leftarrow \text{MAP}(p_{peer}^{(i)})$  ▷ Fix  $p_{peer}^{(i)}$  and optimize  $\{p_{ext}\}_t$  for every window.
8:      $\Delta_{peer}^{(i)} \leftarrow p_{peer}^{(i)} - p_{peer}^{(i-1)}$ 
9:      $\Delta_{ext}^{(i)} \leftarrow \sum_{t=1}^T (p_{ext}^{(i)}(t) - p_{ext}^{(i-1)}(t))$ 
10:     $i \leftarrow i + 1$ 
11:  end while
12:  return  $p_{peer}^{(i)}, \{p_{ext}\}_t$  ▷ The parameters of endogenous and exogenous influence.
13: end procedure

```

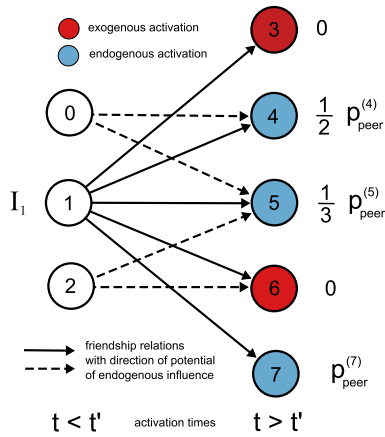


FIGURE 4. Individual and collective influence. In this simple example we estimate influence I_1 of user $i = 1$ as the extent to which he is responsible for endogenous activation $p_{peer}^{(i)}$ of all of his peers $j = \{3, 4, 5, 6, 7\}$ which activated after him. Only three of his peers $j = \{4, 5, 7\}$ activated due to endogenous influence, but he has to share part of this claim with two users $i = \{0, 2\}$ which are their shared peers. The total individual influence for user $i = 1$ in the above example is $I_1 = 1/2 p_{peer}^{(4)} + 1/3 p_{peer}^{(5)} + p_{peer}^{(7)}$. Type of activation (endogenous or exogenous) for each user can be estimated with our methodology or taken from raw data by using referral links from which users visited our application, in which case $p_{peer}^{(i)}$ simply takes values 0 or 1.

just a potential for responsibility and so the user has to share part of his claim to $I^{(i \rightarrow j)}$ with all other m peers of j . For the SI model we can set this to 1, meaning that we consider all peers equally responsible regardless of the time of their activation. Each user would then be assigned $1/m$ of the peer activation probability $p_{peer}^{(j)}$ for each of his peers that activated before him, where m is the number of user's j peers that activated before him. For the EXP model we can weight this with the times of activation - users can claim larger part of the influence for peers that activated close in time to their own activation (more details in Section S4 of the Supporting information). The collective influence for a group of users G is just an average influence of all users in the group $1/|G| \sum_{i \in G} I^{(i)}$.

E. EVALUATION

Instead of using a single threshold for the exogenous responsibility to classify users into endogenously and exogenously

activated we calculate the entire receiver operating characteristic (ROC) curve and associated area under the curve (AUC) score. This allows us to compare different endogenous influence models regardless of the chosen threshold. In order to calculate the ROC curve and AUC score we also need some sort of a gold standard label for each user, for which we use referral links available for sabor2015 and sabor2016 datasets. Depending on the referral link we classify users in one of the three categories: (i) strong endogenous influence for users whose referral link originates from a Facebook share, (ii) potential endogenous influence for users whose referral link originates from Facebook and (iii) strong exogenous influence for users whose referral link originates from an external web site. Users who do not have a referral link are considered as unknown. For the purpose of evaluation we consider users from category (i) as endogenously activated and users from category (iii) as exogenously activated.

F. DATA COLLECTION

Our online survey applications were actually web applications which used Facebook Graph API [47] for authentication of users. Some sort of user authentication was necessary to prevent multiple voting. Facebook Graph API allowed us to collect Facebook friendship relationship between users registered on our application. In addition, with referendum2013.hr we collected basic demographics information such as age and gender. With other three applications we used our own web server directly to collect referral links through which users visited our web application. In all applications we also collect exact registration times of all users. Users provided informed consent on two levels. First, initial web page of survey application displayed a disclaimer next to the registration button describing the type of data we collect and the purpose we intend to use it. This was visible to both registered and unregistered users, before any data was actually collected. Second, upon authorization with their Facebook credentials, but before any data was collected, users were presented with a link to both Facebook's Platform Policy and our own privacy policy, and were given an option to opt out from the survey. In addition, we also provided detailed terms

of use, Frequently Asked Questions (FAQ) and privacy policy web pages which complied with the Facebook's Platform Policy [48], all of which were available to both registered and unregistered users. More details on the data we collected and the methodology of eliciting informed consent is available in Section S2 of the Supporting information. Facebook's Graphs API assigns application-specific ID's to each user, so it is not possible to associate users from different datasets. After they registered users were able to see summary voting statistics of their friends as well as for all registered users. These statistics were displayed *after* the user cast his vote in order to minimize the influence on his choice. We also provided an additional incentive to share the link to the application through Facebook and other social media by displaying to each user a number of users which registered to the application after following the referral link from their share, and comparing this to other users. Our data collection procedure complies with the Facebook Platform Policy [48] and was approved by the Ethics committee of the Faculty of Electrical Engineering and Computing, University of Zagreb.

G. CODE AND DATA AVAILABILITY

Due to Facebook's Platform Policy <https://developers.facebook.com/policy> regarding user's data privacy we are not allowed to publicly release any Facebook-derived data, including personal information and friendship relations between our users. Friendship networks, registration times and analysis code needed to reproduce the results of this paper are available upon a request after signing the data access agreement on <https://goo.gl/forms/IxINFkeBSJpDuzRv2>. The agreement states that the requester: (i) Will only use the dataset for the purpose of reproducing and validating the results of our study; (ii) Will not attempt to deanonymize the dataset or in any other way compromise the identity or privacy of users contained in it; and (iii) Will not further share, distribute, publish, or otherwise disseminate the dataset. This data access agreement complies with the Facebook Platform Policy. Facebook online survey applications through which we collected referendum2013 and sabor2015 datasets are available on public Github repositories: <https://github.com/devArena/referendum2013.hr>, and <https://bitbucket.org/marin/sabor2015.hr>. More information is available in Sections S1 and S2 of the Supporting information.

IV. RESULTS AND DISCUSSION

A. MAXIMUM LIKELIHOOD INFERENCE FOR ENDOGENOUS AND EXOGENOUS INFLUENCE

We want to compute a single set of endogenous influence parameters for the whole period and a separate set of exogenous influence parameters for every time window. Our assumption is that endogenous influence parameters do not change over time, but that exogenous do. A direct way to do this is to perform a joint optimization of a log-likelihood that contains a single set of endogenous influence parameters and a separate set of exogenous influence parameters

for each time window $[t + \Delta t]$. Our log-likelihood would then be $t + 1$ -dimensional in the case of SI model, and $t + 2$ -dimensional for the EXP model - t parameters of exogenous influence for each time window we are considering in our inference plus the parameters of endogenous influence (p_0 for SI model and (p_0, λ) for EXP model). This makes the number of parameters proportional to the number of time windows, which makes a joint optimization of log-likelihood unfeasible. Instead, we use an alternating method [39] described in Algorithm 1 where we alternatively fix either endogenous influence parameters or exogenous influence parameters and optimize the other until both values converge. In addition, we never optimize all of the t parameters of the exogenous influence jointly but do it one by one. This yields a nonparametric estimate for exogenous influence, meaning that we have a separate estimate of exogenous influence $p_{ext}(t)$ at each time step t . Although the number of parameters we have to infer is still proportional to the number of time windows we are considering in our inference, this strategy is much more efficient than joint inference and provides reliable estimates even though there is no formal guarantee that the estimates will actually converge. However, in our experiments we did not experience any problems with the convergence. Fig 5 shows the initialization step of the alternating procedure on a simple simulated activation cascade, where parameters for endogenous and exogenous influence are inferred separately for each time step t .

Using efficient optimization routines allows our method to scale to networks of over 10000 users with resolution of 100 time steps. In our experiments we use a truncated Newton algorithm [45] for maximum likelihood estimation, although in principle any suitable optimization algorithm could be used (more details in Methods section). Total number of users activated due to endogenous and exogenous influence is calculated through the *exogenous responsibility* measure (Eq 4) which is derived from the inferred parameters and quantifies the extent to which is each user's activation is due to endogenous or exogenous influence. This estimate is normalized with the total number of user activations in a given time interval, which is an observable quantity.

B. INFERENCE OF ENDOGENOUS AND EXOGENOUS INFLUENCE ON SIMULATED DATA

Our simulations are designed to approximate, as well as possible, the conditions in which real data were collected. However, instead of using one of the empirical social networks which we collected, we decided to simulate on a configuration model of referendum2013 Facebook friendship network so that our results are reproducible using only a degree sequence. This technique is similar to generating a synthetic representation of a Facebook social network [49] with respect to compactness and anonymity. Configuration model of a network preserves the number of connections each user has, but these connections are permuted randomly across all users. This destroys mesoscale structures such as communities, but is still preferable to other permutation methods where

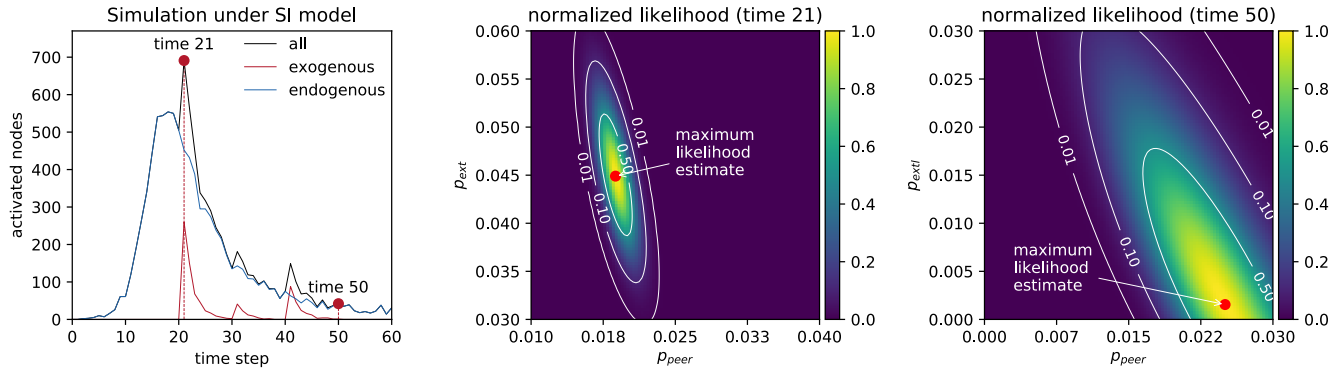


FIGURE 5. Maximum likelihood inference of endogenous and exogenous influence. Plots of the normalized likelihood function (similar to Eq 3 which shows log-likelihood) at two distinct time steps in the simulated activation cascade using SI model for endogenous influence. SI model features only two parameters at each time step - parameter of endogenous influence p_{peer} (p_0 in Eq 1) and a parameter of exogenous influence p_{ext} . Shape of the likelihood function suggests that these two parameters are correlated as each provides part of the explanation for the observed data, and if one is weaker the other most compensate. Also, when we have more data (time 21) the shape of the log-likelihood function is more concentrated than when we have less (time 50), resulting in more confident estimates. In this simulation we are estimating parameters of endogenous and exogenous influence at each time step separately, which corresponds to the initialization stage of our actual inference procedure which we use on simulated (Section IV-B) and empirical (Section IV-C) data. In our full inference procedure we infer a single set of endogenous influence parameters for the whole observation period instead of having a separate estimate for each time step like in this example (more details in Section III-B). Here we are using a truncated Newton algorithm [45] for optimizing a log-likelihood function in order to obtain a maximum likelihood solution, although in practice any suitable optimization method could be used.

either times of activation are permuted (destroying order of activity) or connections themselves are permuted between the users (destroying degree distribution by changing it to binomial) [50]. The simulation starts with a small number of active users and progresses in discrete steps following one of the endogenous influence models (Eq 1 and Eq 2). Fig 6 shows the results using the EXP model (Eq 2) for endogenous influence. At three distinct times we also simulate an exponentially decaying exogenous influence which acts equally on all inactive users. This resembles a typical situation when a distinct exogenous information source activates some of the users [51], which we also observe in our dataset (Fig 2). However, our methodology works equally well for other shapes of exogenous influence (Fig S8 and Fig S9 in the Supporting information). Using just the activation times of all users and their friendship network we are able to estimate the parameters of the assumed endogenous and exogenous influence models as well as the absolute number of users activated predominantly due to the one or the other. In addition, using a measure of *exogenous responsibility* (Eq 4) we are able to infer, for user, the extent to which endogenous or exogenous influence was responsible for activation. Instead of using a single threshold to classify users we calculated the whole receiver operating characteristic (ROC) curve and the corresponding area under the curve (AUC) score to evaluate the performance (Fig 6). We compare our method to a simple baseline commonly used in previous work [39], [52] where an activation is considered exogenous if activated user had no other active peers at the time of the activation. However, as more and more users becomes active, it becomes increasingly likely that a user is connected with at least one other active user by pure chance. This underestimates the number of users activated by exogenous influence and consequently underestimates overall exogenous influence. We obtain similar results (Fig S6 in the Supporting information) for the

SI endogenous influence model and an additional threshold model we define in the Eq S11 of the Supporting information. The inference itself is fast and scales well to networks of over ten thousand users (Section S5 in the Supporting information).

C. INFERENCE OF ENDOGENOUS AND EXOGENOUS INFLUENCE ON EMPIRICAL DATASETS

In order to investigate social interactions between users of a large online social network we developed three online surveys that use Facebook API for collection of data. Surveys were related to three distinct political events in Croatia: 1) *referendum2013.hr* for referendum on definition of marriage, 2) *sabor2015.hr* for parliamentary elections in 2015, and 3) *sabor2016.hr* for parliamentary elections in 2016. Fig 1 and Fig 2 show the collected friendship networks between Facebook users and the number of registrations in 30-minute intervals for each of the survey applications during a week preceding the actual elections. Table 1 shows summary statistics for each of the datasets. The referral links provide information whether each user followed a link originating from a post on Facebook which indicates endogenous influence, or some external website reporting on our survey which indicates exogenous influence. We use this information to evaluate our estimates of endogenous and exogenous influence acting on users. More details on the datasets and the methodology of data collection is available in Section III-B and Sections S1 and S2 of the Supporting information.

Fig 7 shows the results of applying our inference methodology to estimate the magnitude of endogenous and exogenous influence during these three activation cascades. In this experiment we use the EXP model as endogenous influence model because it performed best on average over all three empirical datasets, with and without correction for the observer bias.

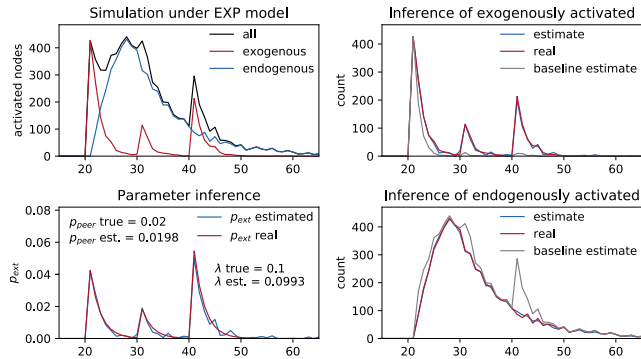


FIGURE 6. Inference on a simulated activation cascade. We use our methodology to infer which users activated due to endogenous or exogenous influence in a simulated activation cascade following exponential decay (EXP) endogenous influence model. In real world applications only total number of activated users (black line) is actually observed, along with the friendship network between users (Fig 6). We use a configuration model of referendum2013 social network to make our results reproducible even without the whole empirical network. We see that our measure is able to differentiate absolute numbers of endogenously and exogenously activated users throughout the whole cascade period and to correctly infer the parameters of endogenous influence - p_{peer} and λ , and exogenous influence $p_{ext}(t)$ for every time period t . We also infer activation type for each user individually by using the exogenous responsibility measure $R^{(i)}(t)$ (Eq 4) as shown on Fig 6 and achieve AUC of 0.93. We compare this with the baseline method where, instead of exogenous responsibility, we use number of active peers at the time of activation. A special case of this baseline is where we consider users without any active peers as exogenously activated, which is a baseline that we use in Fig 6. This baseline method underestimates the exogenously activated users towards the end of the observation period, which is due to the fact that more and more users are active and it is increasingly likely that at least one of the peers is active by chance alone. On Fig 6 we show a histogram of the number of active peers and compare it with exogenous responsibility to demonstrate that no reasonable threshold could not serve as a classification measure, which is also confirmed with a relatively low AUC score of 0.86. The results for SI endogenous influence model are similar and are available in Fig S6 in the Supporting information.

The results for other models are included in Fig S12 and Fig S13 of the Supporting information. As our methodology operates in discrete time (Eq 3) we discretized the activation times of users into 30 minutes time intervals to determine which users were active or inactive during each specific interval. Considering the duration of the data collection for each of the surveys, this corresponds to 333 time intervals for referendum2013 dataset, 327 intervals for sabor2015 dataset and 328 intervals for sabor2016 dataset. Each user that registered on one of the online survey application using his Facebook credentials is considered *activated* in the given time period. The referral link from which we visited the website of the survey application will be used as a proxy of endogenous and exogenous influence - referral links from Facebook are considered as endogenous and those from external websites as exogenous. We later use this information for evaluation of our methodology.

We estimate magnitudes of endogenous and exogenous influence and characterize each user as being endogenously or exogenously activated. We use the AUC score to evaluate the predictive performance of our inferred model on sabor2015 and sabor2016 datasets for which we had data on referral links from which users visited our

survey application. This served as a proxy for ground truth labels which we needed for calculating the AUC scores. The purpose of the model is to estimate the magnitude of endogenous and exogenous influence on each given user, given available data and provided that underlying assumptions of our statistical methodology are satisfied. Similar as in simulated experiments, we compare our methodology with a baseline method that simply estimates the number of exogenously activated users as all those who did not have any active peers at the time of their own activation, and again we observe that it underestimates the number of exogenously activated users, especially near the end of the observation period. Our estimates of endogenously activated users (Fig 7) closely resemble the true number of users activated by following another user's share, which is the strongest indication of endogenous influence we have. On the other hand, it might seem that our method overestimates exogenously activated users by declaring many of the users originating from Facebook as exogenously activated. However, relying on Facebook referrals alone is not a reliable proxy for endogenous activation, as many users might be activated through other means of indirect communication available through Facebook - by following an advertisement, or by directly visiting a Facebook page of the survey application.

We observe that the magnitude of exogenous influence increases as we approach the end of the activation cascade period. This effect is due to the fact that we only observe the friendship network of users that eventually registered on our application, which is only a small subset of the whole Facebook network. However, one of our assumption is that exogenous influence acts uniformly on all users in the friendship network, not just the subset of them, and this manifests in the increased exogenous influence as the activation cascade approaches the size of the network. This *observer bias* can be corrected by adding a correction factor c to our log-likelihood function (Eq 3), which is regulated with parameter α . The results of applying the correction term on the empirical data are shown on Fig 7, while more detailed experiments are available in Fig S5 of the Supporting information). However, because less and less users got activated near the end of the observation period this observer bias does not influence our final estimates by much. However, we still believe that correction is warranted and useful, especially for estimates near the end of the observation period, and in other use cases where observation period is shorter and observer bias might be more pronounced.

For evaluation (Fig 7) we again calculate the corresponding AUC score which uses exogenous responsibility measure $R^{(i)}(t)$ (Eq 4) to classify users into endogenously and exogenously activated. The achieved AUC scores for our method (AUC_{our}) for sabor2015 and sabor2016 datasets are 0.76 and 0.82 respectively. This is higher than the baseline measure which uses number of active peers at the time of activation which achieves AUC scores (AUC_{base}) of 0.68 and 0.78 for the sabor2015 and sabor2016 datasets respectively. Using exponential decay model for endogenous influence allows us

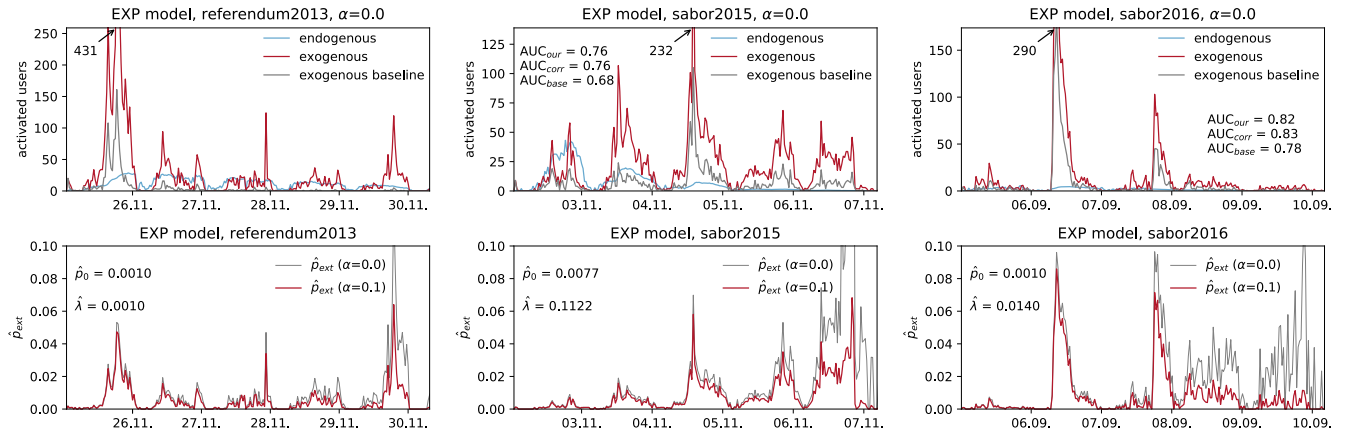


FIGURE 7. Inference on Facebook activation cascades with EXP model. Inference of endogenous and exogenous influence on activation cascades derived from referendum2013, sabor2015 and sabor2016 online survey applications, with EXP model as assumed endogenous influence model. The results for the SI endogenous influence model are in Fig S12 and Fig S13 of the Supporting information. On the bottom panels we see the effect of correction for the observer bias ($\alpha = 0.1$) as compared to no correction ($\alpha = 0$) - it reduces the overestimate of exogenous influence near the end of the observation period. AUC scores for using exogenous responsibility as a measure for classifying users into endogenously and exogenously activated (AUC_{our}) for datasets where we have information on referral links for evaluation - sabor2015 and sabor2016, are 0.76 and 0.82 respectively. This is higher than those achieved with a baseline measure of number of active friends, which are 0.68 and 0.78 for sabor2015 and sabor2016 datasets respectively. A more direct comparison with the baseline is available in Fig S14 of the Supporting information. Facebook referrals alone are not discriminating enough as there are multiple possible ways by which Facebook users might reach our application, including visiting the webpage of our application directly or through an advertisement, both which are more similar to exogenous rather than endogenous influence.

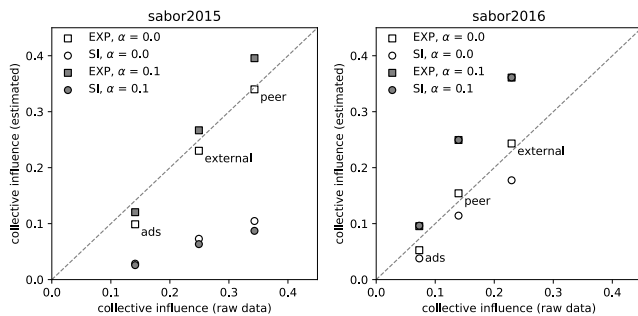


FIGURE 8. Comparison of influence estimates. Comparison of influence estimates obtained from our methodology and raw data for different groups of users - those activated due to endogenous (peer) influence, exogenous (external) influence and advertisements (ads). Ads are similar to exogenous influence as they are targeting large number of users independent of their friendship connections, but within the Facebook social network itself.

to calculate the half-decay of endogenous influence which is 10.1 hours for the sabor2015 dataset. This value is consistent with what we could expect, as it means that endogenous influence diminishes to a fraction of a value in the span of a day or two and requires influx of new users to keep it sustained.

D. COLLECTIVE INFLUENCE

Once we characterized activation of each user as being endogenously or exogenously driven, we can estimate the extent to which each user contributed to the activation of its peers by excluding the portion of the influence attributed to exogenous factors. We do not have a deterministic propagation path for our activation cascade - we do not know who influenced whom directly, so we cannot deterministically incorporate influence of all users in a transitive manner [46].

Nevertheless, our measure of influence simply incorporates all *possible* endogenous propagation paths to estimate an influence for each user (Fig 4 and Eq 5). If we then average this influence over a group of users we get their *collective influence*. Instead of using our estimates of endogenous and exogenous activation for each user we could also estimate influence directly from data by using the referral links from which users visited our application. Fig 8 shows the comparison of our methodology with estimates of influence obtained from raw data for different groups of users that activated due to: endogenous factors, exogenous factors, advertisements. Our question was: Which channel of communication is the most influential, that is, recruits users with higher collective influence? The results of our experiments on two datasets for which we had data on referral links, shows no clear pattern of influence. Different groups of users are more influential depending on the dataset. However, regardless of the model of endogenous influence (SI or EXP) our estimates are robust and are proportional to the ones obtained from raw data. It is important to emphasize again that our methodology does not use any information on referral links or external influence whatsoever, but rather infers this from the dynamics of the user activations. More details is available in Section S4 of the Supporting information.

V. CONCLUSION

Unlike traditional survey methods where data is manually entered either by a respondent or experimenter [53], online social networks provide an opportunity to collect much larger amounts of data on user activity. However, due to their nature they provide challenges to experimental design [54]. Observational studies without explicit consent are regularly performed within companies for marketing purposes, which

is regulated by company's privacy policy, and in some cases this research can be used for academic purposes [36]. Still, academic publication of such research could raise ethical concerns [2], [55]. On the other hand, conducting a study where explicit consent is mandatory heavily restricts the amount of data that can be collected, even when researchers have a direct access to the whole online social network and are in position to present their experiment automatically to the large number of users. For example, a study from Aral and Walker [56] on a sample of 1.3 million Facebook users managed to collect responses of only 7730 users. However, major publicized events such as elections and referendums can serve as catalyzers for mobilizing users. Users are usually willing to participate in a study if through it they receive an information or a service which they perceive as valuable and which could not be easily obtained in some other way.

Despite inherent difficulties in collecting data, we decided to conduct several online surveys using our own web applications and Facebook's API, which allowed us to collect activation cascades and friendship connections of over 20 thousand users in total. Although computational social science is in its infancy, with standards and practices still taking shape, we tried to keep the privacy of the users and follow current recommended ethical practices [9], [10]. Conducting a survey through an online social network means that the recruitment happens organically from person to person as a form of snowball sampling and not through some unbiased randomized procedure, so it's the most eager persons that are recruited first. Number of mobilized users mostly depends on highly connected and willing individuals, that mobilize less willing users. This effect might easily dominate the one from mass media [57].

Using this data we demonstrate how to estimate exogenous and endogenous influence using only information on the friendship connections between users and a single activation cascade which corresponds to the times of user registration. Our methodology exploits the different ways of how exogenous and endogenous influence propagate - endogenous influence propagates between users and as such is dependent on the friendship structure, while exogenous influence acts uniformly on all users regardless of the social network structure. Our method is not able to reconstruct an exact propagation pathway, as these inevitably include pathways external to the particular online social network as well as pathways that are inherently unobservable such as word-of-mouth communication. Still, our method is able to give a probabilistic estimates of these two influences given minimal assumptions. Any additional information on the activation cascade or the social network could be included in our methodology, most probably along the lines of the unified model of social influence [22]. The advantage of such likelihood-based approaches is that inference is performed in a probabilistically-consistent manner, instead of relying on aggregated statistics to choose among competing models of influence [58]. The availability of efficient numerical solvers means our method can easily scale to large

networks of over 10000 users. Computational scalability was already addressed for the unified model [59], however, only for the modeling and not for inference. Our methodology could be applied for characterizing the types of influence in information spreading, for example the role of external factors in the fake news spreading occurring over online social networks such as Facebook or Twitter [60]. Also, there might also be applications outside the domain of social networks as the paradigm of endogenous and exogenous effects could be applied in the wider context of dynamical systems modeling [38].

Our methodology suffers from several limitations, which also indicate potential paths for future research. First, we do not elucidate the mechanisms by which endogenous and exogenous influence arise. The form of the endogenous influence is predefined, and choosing between several possible candidates is possible. In our case, we evaluate different endogenous influence models by their prediction on empirical data, but other methods are possible, including information-theoretic approaches. Second, we assume exogenous influence acts equally on all users, and that parameters of endogenous influence are equal for all users. This was necessary in our case because we only have one activation cascade available for inference [40], and without imposing additional constraints our statistical inference would be infeasible [61], [62]. In cases where multiple activation cascades are available, it should be possible to relax these assumptions and allow for different values of endogenous and exogenous influence parameters for various groups of users. Third, we do not try to correct for the confounding effect arising from unobserved or observed characteristics of users. For example, it is expected that users respond differently to influences, both exogenous and endogenous, from entities that share their political orientation as compared to those that do not. Again, including additional parameters in our model would increase the uncertainty of our estimates. Fourth, we assume friendship connections do not change during the activation cascade. In our case this is justified as the duration of our information cascade is relatively short, only a week, during which we do not expect many changes in friendship connections. For longer observation periods it might be necessary to introduce a possibility of changing friendship connections, which can be incorporated into our model by introducing time-changing quantity $N^{(i)}(t)$ - a set of peers of user i in a particular time step t , in equations for endogenous influence (Eq 1 and 2).

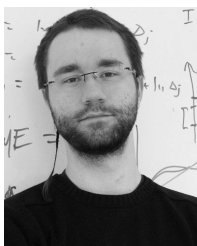
ACKNOWLEDGMENT

We would like to thank the people with whom we had fruitful discussions and who helped in various stages of manuscript preparation and experimental design: Nino Antulov-Fantulin, Vinko Zlatić and Sebastian Krausse. Also, we greatly appreciate the effort of people who actively collaborated in the development of the Facebook online survey applications with which we collected the data: Bruno Rahle, Iva Miholić, Tomislav Lipić, Vedran Ivanac, Matej Mihelčić and Mladen Marinović.

REFERENCES

- [1] J. Borge-Holthoefer, A. Rivero, I. García, E. Cauhé, A. Ferrer, D. Ferrer, D. Francos, D. Iñiguez, M. P. Pérez, G. Ruiz, F. Sanz, F. Serrano, C. Viñas, A. Tarancón, and Y. Moreno, "Structural and dynamical patterns on online social networks: The Spanish May 15th movement as a case study," *PLoS ONE*, vol. 6, no. 8, 2011, Art. no. e23883.
- [2] A. D. I. Kramer, J. E. Guillory, and J. T. Hancock, "Experimental evidence of massive-scale emotional contagion through social networks," *Proc. Nat. Acad. Sci. USA*, vol. 111, no. 24, pp. 8788–8790, 2014.
- [3] K. Lewis, J. Kaufman, M. Gonzalez, A. Wimmer, and N. Christakis, "Tastes, ties, and time: A new social network dataset using Facebook.com," *Social Netw.*, vol. 30, no. 4, pp. 330–342, 2008.
- [4] M. Karsai, G. Iñiguez, K. Kaski, and J. Kertész, "Complex contagion process in spreading of online innovation," *J. Roy. Soc. Interface*, vol. 11, no. 24, pp. 8788–8790, 2014. [Online]. Available: <http://rsif.royalsocietypublishing.org/content/11/101/20140694.abstract>
- [5] A. Guille, H. Hacid, C. Favre, and D. A. Zighed, "Information diffusion in online social networks: A survey," *ACM SIGMOD Rec.*, vol. 42, no. 2, pp. 17–28, 2013.
- [6] M. De Domenico, A. Lima, P. Mougél, and M. Musolesi, "The anatomy of a scientific rumor," *Sci. Rep.*, vol. 3, Oct. 2013, Art. no. 2980.
- [7] A. Najar, L. Denoyer, and P. Gallinari, "Predicting information diffusion on social networks with partial knowledge," in *Proc. 21st Annu. Conf. World Wide Web Companion (WWW)*, 2012, pp. 1197–1204.
- [8] J. Yang and J. Leskovec, "Modeling information diffusion in implicit networks," in *Proc. IEEE Int. Conf. Data Mining (ICDM)*, Dec. 2010, pp. 599–608.
- [9] M. J. Salganik, *Bit by Bit: Social Research in the Digital Age*. Princeton, NJ, USA: Princeton Univ. Press, 2017.
- [10] M. Kosinski, S. C. Matz, S. D. Gosling, V. Popov, and D. Stillwell, "Facebook as a research tool for the social sciences: Opportunities, challenges, ethical considerations, and practical guidelines," *Amer. Psychol.*, vol. 70, no. 6, pp. 543–556, 2015.
- [11] M. Granovetter, "Threshold models of collective behavior," *Amer. J. Sociol.*, vol. 83, no. 6, pp. 1420–1443, 1978.
- [12] D. J. Watts, "A simple model of global cascades on random networks," *Proc. Nat. Acad. Sci. USA*, vol. 99, no. 9, pp. 5766–5771, 2002.
- [13] S. Aral and C. Nicolaides, "Exercise contagion in a global social network," *Nature Commun.*, vol. 8, Apr. 2017, Art. no. 14753.
- [14] J. Leskovec, L. A. Adamic, and B. A. Huberman, "The dynamics of viral marketing," *ACM Trans. Web*, vol. 1, no. 1, 2007, Art. no. 5.
- [15] M. Kimura, K. Saito, K. Ohara, and H. Motoda, "Learning information diffusion model in a social network for predicting influence of nodes," *Intell. Data Anal.*, vol. 15, no. 4, pp. 633–652, 2011.
- [16] D. Centola, "The spread of behavior in an online social network experiment," *Science*, vol. 329, no. 5996, pp. 1194–1197, 2010.
- [17] R. Pastor-Satorras and A. Vespignani, "Epidemic spreading in scale-free networks," *Phys. Rev. Lett.*, vol. 86, no. 14, pp. 3200–3203, Apr. 2001.
- [18] S. Aral, L. Muchnik, and A. Sundararajan, "Distinguishing influence-based contagion from homophily-driven diffusion in dynamic networks," *Proc. Nat. Acad. Sci. USA*, vol. 106, no. 51, pp. 21544–21549, 2009.
- [19] S. González-Bailón, J. Borge-Holthoefer, A. Rivero, and Y. Moreno, "The dynamics of protest recruitment through an online network," *Sci. Rep.*, vol. 1, Dec. 2011, Art. no. 197.
- [20] X. Lu and C. Brelsford, "Network structure and community evolution on Twitter: Human behavior change in response to the 2011 Japanese earthquake and tsunami," *Sci. Rep.*, vol. 4, Oct. 2014, Art. no. 6773.
- [21] D. M. J. Lazer, M. A. Baum, Y. Benkler, A. J. Berinsky, K. M. Greenhill, F. Menczer, M. J. Metzger, B. Nyhan, G. Pennycook, D. Rothschild, M. Schudson, S. A. Sloman, C. R. Sunstein, E. A. Thorson, D. J. Watts, and J. L. Zittrain, "The science of fake news," *Science*, vol. 359, no. 6380, pp. 1094–1096, 2018.
- [22] A. Srivastava, C. Chelms, and V. K. Prasanna, "The unified model of social influence and its application in influence maximization," *Social Netw. Anal. Mining*, vol. 5, no. 1, p. 66, 2015.
- [23] D. Kempe, J. Kleinberg, and E. Tardos, "Maximizing the spread of influence through a social network," in *Proc. 9th ACM SIGKDD Int. Conf. Knowl. Discovery Data Mining*, 2003, pp. 137–146.
- [24] M. Piškorec, N. Antulov-Fantulin, I. Miholic, T. Šmuc, and M. Šikić, "Modeling peer and external influence in online social networks: Case of 2013 referendum in Croatia," in *Complex Networks & Their Applications VI*, C. Cherifi, H. Cherifi, M. Karsai, and M. Musolesi, Eds. Cham, Switzerland: Springer, 2018, pp. 1015–1027.
- [25] J.-P. Onnela and F. Reed-Tsochas, "Spontaneous emergence of social influence in online systems," *Proc. Nat. Acad. Sci. USA*, vol. 107, no. 43, pp. 18375–18380, 2010.
- [26] A. Goyal, F. Bonchi, and L. V. S. Lakshmanan, "Learning influence probabilities in social networks," in *Proc. 3rd ACM Int. Conf. Web Search Data Mining (WSDM)*, 2010, pp. 241–250.
- [27] D. J. Daley and D. G. Kendall, "Stochastic rumors," *IMA J. Appl. Math.*, vol. 1, no. 1, pp. 42–55, 1965.
- [28] D. P. Maki-Thompson, *Mathematical Models and Applications, With Emphasis on Social, Life, and Management Sciences*. Upper Saddle River, NJ, USA: Prentice-Hall, 1973.
- [29] A. L. Hill, D. G. Rand, M. A. Nowak, and N. A. Christakis, "Infectious disease modeling of social contagion in networks," *PLoS Comput. Biol.*, vol. 6, no. 11, 2010, Art. no. e1000968.
- [30] J. Goldenberg, B. Libai, and E. Müller, "Talk of the network: A complex systems look at the underlying process of word-of-mouth," *Marketing Lett.*, vol. 12, no. 3, pp. 211–223, 2001.
- [31] H. Narasimhan, D. C. Parkes, and Y. Singer, "Learnability of influence in networks," in *Proc. Adv. Neural Inf. Process. Syst.*, C. Cortes, N. D. Lawrence, D. D. Lee, M. Sugiyama, and R. Garnett, Eds. Red Hook, NY, USA: Curran Associates, 2015, pp. 3186–3194. [Online]. Available: <http://papers.nips.cc/paper/5989-learnability-of-influence-in-networks.pdf>
- [32] C. Wang, W. Chen, and Y. Wang, "Scalable influence maximization for independent cascade model in large-scale social networks," *Data Mining Knowl. Discovery*, vol. 25, no. 3, pp. 545–576, Nov. 2012.
- [33] J. Yang and J. Liu, "Influence maximization-cost minimization in social networks based on a multiobjective discrete particle swarm optimization algorithm," *IEEE Access*, vol. 6, pp. 2320–2329, 2017.
- [34] D. Guillebaud, J. Becker, and D. Centola, "Complex contagions: A decade in review," in *Complex Spreading Phenomena in Social Systems*, S. Lehmann and Y.-Y. Ahn, Ed. Cham, Switzerland: Springer, 2018, pp. 3–25.
- [35] A. Mahmoodi, B. Bahrami, and C. Mehring, "Reciprocity of social influence," *Nature Commun.*, vol. 9, no. 1, Jun. 2018, Art. no. 2474.
- [36] D. Eckles, R. F. Kizilcec, and E. Bakshy, "Estimating peer effects in networks with peer encouragement designs," *Proc. Nat. Acad. Sci. USA*, vol. 113, no. 27, pp. 7316–7322, 2016.
- [37] C. R. Shalizi and A. C. Thomas, "Homophily and contagion are generically confounded in observational social network studies," *Sociol. Methods Res.*, vol. 40, no. 2, pp. 211–239, 2011.
- [38] M. A. de Menezes and A.-L. Barabási, "Separating internal and external dynamics of complex systems," *Phys. Rev. Lett.*, vol. 93, no. 6, pp. 068701-1–068701-4, 2004.
- [39] S. A. Myers, C. Zhu, and J. Leskovec, "Information diffusion and external influence in networks," in *Proc. 18th ACM SIGKDD Int. Conf. Knowl. Discovery Data Mining*, 2012, pp. 33–41.
- [40] M. Karsai, G. Iñiguez, R. Kikas, K. Kaski, and J. Kertész, "Local cascades induced global contagion: How heterogeneous thresholds, exogenous effects, and unconcerned behaviour govern online adoption spreading," *Sci. Rep.*, vol. 6, Jun. 2016, Art. no. 27178.
- [41] A. Anagnostopoulos, G. Brova, and E. Terzi, "Peer and authority pressure in information-propagation models," in *Machine Learning and Knowledge Discovery in Databases (Lecture Notes in Computer Science: Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, vol. 6911. Berlin, Germany: Springer, 2011, pp. 76–91.
- [42] P. Brach, A. Epasto, A. Panconesi, and P. Sankowski, "Spreading rumours without the network," in *Proc. 2nd ACM Conf. Online Social Netw. (COSN)*, 2014, pp. 107–118.
- [43] T. Takaguchi, N. Masuda, and P. Holme, "Bursty communication patterns facilitate spreading in a threshold-based epidemic dynamics," *PLoS ONE*, vol. 8, no. 7, 2013, Art. no. e68629.
- [44] F. Karimi and P. Holme, "Threshold model of cascades in empirical temporal networks," *Phys. A, Stat. Mech. Appl.*, vol. 392, no. 16, pp. 3476–3483, 2013.
- [45] S. G. Nash, "A survey of truncated-Newton methods," *J. Comput. Appl. Math.*, vol. 124, nos. 1–2, pp. 45–59, Dec. 2000.
- [46] X. Teng, S. Pei, F. Morone, and H. A. Makse, "Collective influence of multiple spreaders evaluated by tracing real information flow in large-scale social networks," *Sci. Rep.*, vol. 6, Oct. 2016, Art. no. 36043.
- [47] *Graph API*. Accessed: Jul. 26, 2018. [Online]. Available: <https://developers.facebook.com/docs/graph-api>
- [48] *Facebook Platform Policy*. Accessed: Jul. 26, 2018. [Online]. Available: <https://developers.facebook.com/policy>

- [49] L. Humski, D. Pinter, and M. Vranić, "Analysis of Facebook interaction as basis for synthetic expanded social graph generation," *IEEE Access*, vol. 7, pp. 6622–6636, 2018.
- [50] P. Holme, "Modern temporal network theory: A colloquium," *Eur. Phys. J. B*, vol. 88, no. 9, p. 234, 2015.
- [51] R. Crane and D. Sornette, "Robust dynamic classes revealed by measuring the response function of a social system," *Proc. Nat. Acad. Sci. USA*, vol. 105, no. 41, pp. 15649–15653, 2008.
- [52] M. Gomez-Rodriguez, D. Balduzzi, and B. Schölkopf, "Uncovering the temporal dynamics of diffusion networks," in *Proc. 28th Int. Conf. Mach. Learn. (ICML)*, 2011, pp. 561–568.
- [53] A. Dhand, C. C. White, C. Johnson, Z. Xia, and P. L. De Jager, "A scalable online tool for quantitative social network assessment reveals potentially modifiable social environmental risks," *Nature Commun.*, vol. 9, no. 1, 2018, Art. no. 3930.
- [54] D. Walker and L. Muchnik, "Design of randomized experiments in networks," *Proc. IEEE*, vol. 102, no. 12, pp. 1940–1951, Dec. 2014.
- [55] I. M. Verma, "Editorial expression of concern: Experimental evidence of massive-scale emotional contagion through social networks," *Proc. Nat. Acad. Sci. USA*, vol. 111, no. 29, 2014, Art. no. 10779.
- [56] S. Aral and D. Walker, "Identifying influential and susceptible members of social networks," *Science*, vol. 337, no. 6092, pp. 337–341, 2012.
- [57] A. Rutherford, M. Cebrian, S. Dsouza, E. Moro, A. Pentland, and I. Rahwan, "Limits of social mobilization," *Proc. Nat. Acad. Sci. USA*, vol. 110, no. 16, pp. 6281–6286, 2013.
- [58] J. Leskovec, L. Backstrom, R. Kumar, and A. Tomkins, "Microscopic evolution of social networks," in *Proc. 14th ACM SIGKDD Int. Conf. Knowl. Discovery Data Mining*, 2008, pp. 462–470.
- [59] A. Popa, M. Frincu, and C. Chelmiss, "A distributed algorithm for the efficient computation of the unified model of social influence on massive datasets," in *Proc. IEEE High Perform. Extreme Comput. Conf. (HPEC)*, Sep. 2017, pp. 1–7.
- [60] A. Bovet and H. A. Makse, "Influence of fake news in Twitter during the 2016 US presidential election," *Nature Commun.*, vol. 10, no. 1, 2019, Art. no. 7.
- [61] Y. Yoshikawa, K. Saito, H. Motoda, K. Ohara, and M. Kimura, "Acquiring expected influence curve from single diffusion sequence," in *Knowledge Management and Acquisition for Smart Systems and Services (Lecture Notes in Computer Science: Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, vol. 6232. Berlin, Germany: Springer, 2010, pp. 273–287.
- [62] N. Du, Y. Liang, M.-F. Balcan, and L. Song, "Influence function learning in information diffusion networks," in *Proc. 31st Int. Conf. Mach. Learn. (ICML)*, vol. 5, 2014, pp. 4118–4135.



MATIJA PIŠKOREC received the M.Sc. degree in computer science from the Faculty of Electrical Engineering and Computing, University of Zagreb, where he is currently pursuing the Ph.D. degree. He is also a Research Associate with the Laboratory for Machine Learning and Knowledge Representation, Rudjer Boskovic Institute, Zagreb, and is a member of the Centre of Excellence Project "DATACROSS". His main research

interests are in the field of machine learning and complex systems, with special emphasis on inference of dynamical processes on social networks. He is also interested in building web services related to data analysis and visualization.



TOMISLAV ŠMUC received the Ph.D. degree. He was a Mentor of a dozen of master's and Ph.D. degrees students with the University of Zagreb, involved in organization of international conferences (ECML-PKDD, Discovery Science) on several occasions. He is currently a the Head of the Laboratory for Machine Learning and Knowledge Representation, at Rudjer Boskovic Institute, Zagreb. He has published more than 100 articles in journals and proceedings of international conferences. His research interests are in the area of artificial intelligence, in development and use of machine learning, and data mining techniques for knowledge discovery in different domains of science and technology, for last twenty years. In this period, he has been participating in, or leading, a number of research projects financed by Croatian, European, and other international funding agencies. He is an Evaluator for several research funding agencies. He also serves as a reviewer for a number of scientific journals in the fields of computer science, computational biology, and interdisciplinary science.



MILE ŠIKIĆ received the Ph.D. degree in computer science from the University of Zagreb, in 2008. For the first seven years of his career, he was a System Integrator, Consultant, and Project Manager on the projects with industry in the fields of computer and mobile networks. In 2009, he became an Assistant Professor of computer science. He is currently a Professor with the Faculty of Electrical Engineering and Computing, University of Zagreb, Croatia. He is also spending

his sabbatical year with the Genome Institute of Singapore. His scientific work is focused on the development of new algorithms and machine learning methods for genome sequence analysis and analysis of dynamics in networks.

...