Topic-aware Social Influence Propagation Models

Nicola Barbieri · Francesco Bonchi · Giuseppe Manco

Received: Jan 04, 2013 / Revised: Mar 20, 2013 / Accepted: Mar 29, 2013

Abstract The study of influence-driven propagations in social networks and its exploitation for *viral marketing* purposes have recently received a large deal of attention However, regardless the fact that users authoritativeness, expertise, trust and influence are evidently topic-dependent, the research on social influence has surprisingly largely overlooked this aspect.

In this article we study social influence from a topic modeling perspective. We introduce novel topic-aware influence-driven propagation models that, as we show in our experiments, are more accurate in describing real-world cascades than the standard (i.e., topic-blind) propagation models studied in the literature. In particular, we first propose simple topic-aware extensions of the well-known Independent Cascade and Linear Threshold models. However, these propagation models have a very large number of parameters which could lead to overfitting.

Therefore we propose a different approach explicitly modeling authoritativeness, influence and relevance under a topic-aware perspective. Instead of considering user-to-user influence, the proposed model focuses on user authoritativeness and interests in a topic, leading to a drastic reduction of the number of parameters of the model. We devise methods to learn the parameters of the models from a dataset of past propagations. Our experimentation confirms the high accuracy of the proposed models and learning schemes.

Nicola Barbieri Yahoo! Research E-mail: barbieri@yahoo-inc.com

Francesco Bonchi Yahoo! Research E-mail: bonchi@yahoo-inc.com

Giuseppe Manco ICAR-CNR E-mail: manco@icar.cnr.it Keywords Social Influence \cdot Topic Modeling \cdot Topic-aware Propagation Model \cdot Viral Marketing

1 Introduction

Social influence and the phenomenon of influence-driven propagations in social networks have received tremendous attention in the last years, fueled by a variety of applications, such as *viral marketing* [8,23], personalized recommendations [26], feed ranking [12], and the analysis of Twitter [29,1] just to name a few. One of the key computational problems in this area is the identification of a set of influential users, which are more likely to produce large influence-driven cascades: these are the users that should be "targeted" by a viral marketing campaign. This problem has received a good deal of attention by the data mining research community in the last decade [4], but quite surprisingly, the *characteristics of the item being the subject of the viral marketing campaign has been left out of the picture*.

Kempe *et al.* [13] formalize the *influence maximization* problem for a generic item: for a given budget k, find k "seed" nodes in the network, such that by activating them we can maximize the expected number of nodes that eventually get activated, according to a chosen *propagation model*, that governs how influence diffuses or propagates through the network. Kempe *et al.* [13] mainly focus on two propagation models – the *Independent Cascade* (IC) and the *Linear Threshold* (LT) models. Following this seminal work, a substantial research effort has been dedicated to develop algorithms for influence maximization under these two propagation models (see Section 2). However, these propagation models suffer various limitations when it comes to model real-world cascades: e.g., the discrete treatment of time and the very large number of parameters. The latter is a serious issue both for efficiency and scalability, but more importantly, for the risk of overfitting.

In this paper we start from the observations that (i) users have different interests, (ii) items have different characteristics and (iii) similar items are likely to interest the same users. Following these observations, we take a topicmodeling perspective to jointly learn items characteristics, users' interests and social influence. This results in new propagation models that experimentally are proven to be more accurate in describing real-world cascades.

More in details the contributions of this article are as follows:

- We extend the classic IC and LT models to be topic-aware. The propagation models we obtain are dubbed *Topic-aware Independent Cascade* (TIC) model and *Topic-aware Linear Threshold* (TLT) model.
- For the problem of influence maximization, we show that the objective function to maximize, i.e., the *expected spread*, remains monotone and submodular for both TIC and TLT models. Thus the simple greedy algorithm provides a $(1 - 1/e - \phi)$ -approximation of the optimal solution.

- We devise an *expectation maximization* (EM) approach for estimating the parameters of the TIC model.
- Starting from a discussion on the limits of the TIC and TLT models, we introduce a new influence propagation model, dubbed AIR (*Authoritativeness-Interest-Relevance*). Instead of considering user-to-user influence, the proposed model focuses on user authoritativeness and interests in a topic, leading to a drastic reduction of the number of parameters of the model, with benefits in terms of reduced risk of over-fitting and reduced learning time.
- We devise a generalized expectation maximization (GEM) approach to learn the parameters that maximize the likelihood for the AIR model.
- Our experiments on real-world social networks show that topic-aware influence propagation models outperform the traditional "topic-blind" IC model in predicting adoption of a specific item, thus in modeling real-world cascades.
- The benefits of keeping in consideration the characteristics of the item being propagated, are confirmed by our experiments on influence maximization: topic-aware methods exhibit a consistent gain over state-of-the-art approach that just considers a generic item, ignoring its characteristics.

Although topic-wise social influence has been studied before, to the best of our knowledge we are the first to study it within the context of *viral marketing* and the *influence maximization* problem, and to propose topic-aware influence propagation models.¹ The collocation of our contribution within the state of the art is discussed in details in the next section. In Section 3 we introduce the TIC and TLT models, while Section 4 is devoted to the AIR model. Section 5 reports our experimental analysis. Section 6 discusses future work and concludes the article.

2 Background and Related Work

In this section we provide the needed background for introducing the contributions of this article, while discussing their collocation within the state of the art.

2.1 Influence maximization

Suppose we are given a social network, that is a directed graph whose nodes are users and arcs represent social relations among the users. Suppose we are also given the estimates of reciprocal influence between individuals connected in the network, that is a weight (or probability) $p_{v,u}$ associated top each arc

 $^{^1\,}$ Note that the present manuscript is an invited extended version of our paper presented at the ICDM 2012 conference with the same title [21].

(v, u). As said in the previous section, a basic computational problem is that of selecting the set of initial users that are more likely to influence the largest number of users in the social network. The first algorithmic treatment of the problem was provided by Domingos and Richardson [8,23], who modeled the diffusion process in terms of Markov random fields, and proposed heuristic solutions to the problem.

Later, Kempe *et al.* [13] studied influence maximization as a discrete optimization problem focusing on two fundamental propagation models, named *Independent Cascade Model* (IC) and *Linear Threshold Model* (LT). In both these models, at a given timestamp, each node is either active (an adopter of the innovation, or a customer which already purchased the product) or inactive, and each node's tendency to become active increases monotonically as more of its neighbors become active. An active node never becomes inactive again. Time unfolds deterministically in discrete steps. As time unfolds, more and more of neighbors of an inactive node u become active, eventually making u become active, and u's decision may in turn trigger further decisions by nodes to which u is connected.

In the IC model each arc (u, v) has an associated probability $p_{v,u}$ that can be considered as the strength of the influence of v over u. When a node v first becomes active, say at time t, it is considered contagious. It has one chance of influencing each inactive neighbor u with probability $p_{v,u}$, independently of the history thus far. If the tentative succeeds, u becomes active at time t + 1.

In the LT model, each node u is influenced by each neighbor v according to a weight $p_{v,u}$, such that the sum of incoming weights to u is no more than 1. At the beginning of the propagation each node u chooses a threshold θ_u uniformly at random from [0, 1]. At any timestamp t, if the total weight from the active neighbors of an inactive node u is at least θ_u , then u becomes active at timestamp t + 1.

In both the models, the process repeats until no new node becomes active. Given a propagation model m (e.g., IC or LT) and a seed set $S \subseteq V$, the expected number of active nodes at the end of the process is denoted by $\sigma_m(S)$. The *influence maximization problem* requires to find the set $S \subseteq V$, |S| = k, such that $\sigma_m(S)$ is maximum.

Kempe *et al.* show that under both the IC and LT propagation models, the problem is **NP**-hard [13]. Kempe *et al.*, however, show that the function $\sigma_m(S)$ is monotone (i.e., $\sigma_m(S) \leq \sigma_m(T)$ whenever $S \subseteq T$) and submodular (i.e., $\sigma_m(S \cup \{w\}) - \sigma_m(S) \geq \sigma_m(T \cup \{w\}) - \sigma_m(T)$ whenever $S \subseteq T$). When equipped with such properties, the simple greedy algorithm that at each iteration greedily extends the set of seeds with the node providing the largest marginal gain, produces a solution with provable approximation guarantee (1 - 1/e) [20].

Though simple, the greedy algorithm is computationally prohibitive, since the step of selecting the node providing the largest marginal gain is $\#\mathbf{P}$ -hard under both the IC and the LT model. In their paper, Kempe *et al.* run Monte Carlo simulations for sufficiently many times to obtain an accurate estimate of the expected spread. In particular, they show that for any $\phi > 0$, there is a $\delta > 0$ such that by using $(1+\delta)$ -approximate values of the expected spread, we obtain a $(1-1/e-\phi)$ -approximation for the influence maximization problem. However, running many propagation simulations is extremely costly on very large real-world social networks. Therefore, following [13], considerable effort has been devoted to develop methods for improving the efficiency of influence maximization [14, 15, 5, 10, 24, 6].

The approaches discussed above, assume a weighted social graph as input and do not address *how* the link influence weights (or probabilities) can be obtained; [25, 28, 9, 30] instead focus on the latter problem and propose specific solutions. Saito *et al.* [25] for example, study how to learn the probabilities for the IC model from a set of past propagations. They formalize this as a likelihood maximization problem and then apply the Expectation Maximization (EM) algorithm to solve it. We will extend this contribution to deal with topic-wise influence in Section 3.1.

Goyal *et al.* [9] also study the problem of learning influence probabilities but under a different model, i.e., an instance of the General Threshold Model. They extend this model by introducing temporal decay, as well as factors such as the influenceability of a specific user, and influence-proneness of a certain action. They also show that their methods can be used to predict *whether* a user will perform an action and *when*, and this prediction has higher accuracy (i.e., it's easier) for users with higher influenceability scores.

2.2 Topic modeling

Probabilistic Topic models [3,27,2] include a suite of techniques which widely used in text analysis. They provide a low-dimensional semantic representation that allows the discovering of global relationships within data by exploiting co-occurrence. The key idea at the basis of topic modeling, is to introduce an hidden variable Z for each co-occurrence of words within a corpus of documents. This hidden variable can range among K states and each topic (i.e., state of the latent variable) represents an abstract interest/pattern and intuitively models the underlying cause for each data observation.

Given a corpus, the assumption behind this family of techniques is that each document may exhibit multiple topics and each word in the document is generated by a particular topic. More specifically, each document is represented as a mixture of topics, where each topic induces a distribution over words of the considered dictionary.

Among the probabilistic approaches for topic modeling, besides mixture models that are widely investigated in the literature [7], *Probabilistic Latent Semantic Analysis* (pLSA) [11] is considered the progenitor of a wide range of recent approaches, which include e.g. the popular *Latent Dirichlet Allocation* (LDA) [3]. A probabilistic topic model specifies a generative process for documents which, at high level, can be summarized as follows:: (i) to generate a new document, we generate a distribution over topics; then, (ii) for each word to be generated we (a) choose a topic by drawing upon the document-specific distribution over topics and finally (b) generate a word from the topic-specific distribution over tokens in the dictionary.

The difference between pLSA and LDA relies on how the document-specific distribution over topics is generated: while in pLSA the topic-mixture weights are directly modelled in the inference process, in LDA this distribution is drawn from a Dirichlet distribution with a corpus specific hyperparameter α . This further level of abstraction make easier the generalization of the model to new (unobserved) documents.

2.3 Topic-aware influence analysis

Regardless the fact that users authoritativeness, expertise, trust and influence are evidently topic-dependent, the research on social influence has surprisingly largely overlooked this aspect. To the best of our knowledge only few papers have looked at social influence from the topics perspective [28,29,17,16].

Tang *et al.* [28] study the problem of learning user-to-user topic-wise influence strength. The input to their problem is the social network and a prior topic distribution for each node, which is given as input and inferred separately. As a consequence, they do not consider the simultaneous learning of topics and topic-wise influence. Further, their main focus is expert finding, and hence they do not propose any propagation model, nor study influence maximization. They instead deal with the efficiency problem by devising a distributed learning algorithm under the Map-Reduce programming model.

A probabilistic model for the joint inference of the topic distribution and topic-wise influence strength has been proposed by Liu *et al.* [17]. Here the input is an heterogenous social network with nodes that are users and documents. The goal is to learn users' interest (topic distribution) and user-to-user influence. Gibbs-Sampling algorithm is used to estimate both topic distribution and influence weights.

Lin *et al.* [16] study the joint modeling of influence and topics, by adopting textual models. According to the generative semantic of the proposed approach, each document is generated by a mixture model on topics. The topic sampling process takes into account document-to-documents non negative weights which models influence (in this case topic-inheritance), while novel aspects of the document are modeled by the evolution component.

Weng *et al.* [29] analyze topic-wise influence in **Twitter** by means of a twostep process. First, topics of interest for each user are extracted by means of LDA and topic-specific relationship networks are constructed. Then, in order to measure the influence of each user, they propose TwitterRank, an extension of the PageRank algorithm taking into account both the topic similarity and the social link structure.

What mentioned before for [28] holds for [17][16][29] too: none of these papers define an influence propagation model nor study the influence maximization problem, as we do in the present article. In conclusion, our work collocates in the intersection of the area of research on influence propagation and maximization with the area of topic-modeling and, to the best of our knowledge, it is the first work proposing topic-aware propagation models and algorithms to learn social influence strength and the topics *jointly*.

3 Simple Topic-aware Propagation Models

As a first step towards topic-aware modeling of social influence, we extend the classic Independent Cascade (IC) and Linear Threshold (LT) models to their topic-aware versions. The two standard (i.e., topic-blind) models were introduced in the previous section.

Topic-aware Independent Cascade Model (TIC). In the topic-aware version of the IC model the user-to-user influence probabilities depend on the topic. Therefore, for each arc $(v, u) \in E$ and each topic $z \in [1, K]$ we are given a probability $p_{v,u}^z$, representing the strength of the influence exerted by user v on user u on topic z. Moreover for each item i that propagates in the network, we have a distribution over the topics, that is for each topic $z \in [1, K]$ we are given $\gamma_i^z = P(Z = z|i)$, with $\sum_{z=1}^K \gamma_i^z = 1$.

In this model a propagation happens like in the IC model: when a node v first becomes active on item i, has one chance of influencing each inactive neighbor u, independently of the history thus far. The tentative succeeds with a probability that is the weighted average of the link probability w.r.t. the topic distribution of the item i:

$$p_{v,u}^{i} = \sum_{z=1}^{K} \gamma_{i}^{z} p_{v,u}^{z}.$$
 (1)

Similarly we can formulate a topic-aware version of the LT model.

Topic-aware Linear Threshold Model (TLT). For each arc $(v, u) \in E$ and each topic $z \in [1, K]$ we are given a weight $p_{v,u}^z$, such that the sum of incoming weights in each node and for each topic is no more than 1. Each node u chooses a threshold θ_u uniformly at random from [0, 1]. At time t, a node u which is not yet active on item i, is submitted to an influence weight

$$W_i^t(u) = \sum_{z=1}^K \sum_{v \in \mathcal{F}_i(u,t)} \gamma_i^z p_{v,u}^z.$$

$$\tag{2}$$

where $\mathcal{F}_i(u, t)$ denotes the set of users that have a link to u and that at time t have already adopted the item i. If $W_i^t(u) \ge \theta_u$, then u will activate on item i at time t + 1.

Observation 1 For both the TIC and TLT models the submodularity of the expected spread $\sigma_m(S)$ is directly inherited from the IC and LT models, respectively. In fact, in both cases only the model parameters are topic-aware, while the overall mechanism of propagation does not change. In particular, given an item i just let $p_{v,u} := \sum_{z=1}^{K} \gamma_i^z p_{v,u}^z$ (Eq. 1) to reduce from TIC to IC and from TLT to LT.

N. Barbieri et al.

notation	description
G = (V, E)	directed social graph
$(v,u) \in E$	an arc from user $v \in V$ to user $u \in V$
$\sigma_m(S)$	expected spread of $S \subseteq V$ under model m
I	the universe of items (index i)
i	index over the item-set
K	number of topics
$z \in [1, K]$	a topic
$p_{v,u}^z$	strength of influence of v on u , on topic z
γ_i^z	topic distribution for item i
ϑ_u^z	topic distribution for user u
\mathbb{D}^{-}	input DB of propagations (propagation log)
$(v, i, t) \in \mathbb{D}$	user v adopts item i at time t
$t_i(v)$	the time at which v adopts item i
$D_i(t)$	$\{v \in V t_i(v) = t\}$
t_i and $\overline{t_i}$	min and max t s.t. $D_i(t) \neq \emptyset$
$-C_i(t)$	$\bigcup_{t' \le t} D_i(t')$
$\mathcal{F}_i(u,t)$	$\{v \in V (v, u) \in E \land v \in C_i(t)\}$
$W_i^t(u)$	total influence for item i on u at time t
p_v^z	authoritativeness of user v in topic z
φ_i^z	relevance of item i in topic z
Å	influence window

Table 1: Some of the notation used

Essentially what Observation 1 states is that the mechanism of propagation does not change between IC, LT and their topic-aware counter-part. What changes is the fact that, while in the classic model there is no distinction among different items, in the topic-aware models different items induce different influence strength over the links. Given a specific item *i* and its topic-distribution γ_i , we just can "re-compute" the links strength as $p_{v,u} := \sum_{z=1}^{K} \gamma_i^z p_{v,u}^z$ and then apply the standard, topic-blind, models.

Hence it holds the following.

Proposition 1 The expected spread $\sigma_m(S)$ remains monotone and submodular for m = TIC or m = TLT.

Proof The proof follows directly from the proofs for IC and LT in [13] and Observation 1.

A direct corollary is that the greedy algorithm provides an $(1 - 1/e - \phi)$ -approximation for the influence maximization problem also under the TIC and TLT propagation models [13].

Next we define an Expectation Maximization (EM) method for learning the parameters of the TIC model.

3.1 Learning topic-aware influence

The problem of learning the parameters of the TIC models takes in input the social graph G = (V, E), a log of past propagations \mathbb{D} , and an integer K. The

propagation log is a relation (User, Item, Time) where a tuple $(u, i, t) \in \mathbb{D}$ indicates that user u adopted item i at time t. We assume that no user adopts the same item more than once. Moreover we assume that the projection of \mathbb{D} on User is contained in the set of nodes V of the social graph G. We let \mathcal{I} denote the universe of items, i.e., the projection of \mathbb{D} on the second column. We also use D_i to denote the propagation trace of i, that is the selection of the tuples of \mathbb{D} where Item = i, while $D_i(t)$ will denote the set of users that adopted i at time t, and $C_i(t) = \bigcup_{t' \leq t} D_i(t')$. Finally we use $\underline{t_i}$ and $\overline{t_i}$ to denote the first and last timestamp of adoption of item i.

The output of the learning problem is the set of all parameters of the TIC propagation model, which we denote Θ : these are γ_i^z and $p_{v,u}^z$ for all $i \in \mathcal{I}$, $(v, u) \in E$, and $z \in [1, K]$.

Assuming that each propagation trace is independent from the others, the likelihood of the data given the model parameters Θ , can be expressed as:

$$\mathcal{L}(\Theta; \mathbb{D}) = \sum_{i \in \mathcal{I}} \log \mathcal{L}(\Theta; D_i).$$
(3)

Saito *et al.* [25] assume that the input propagations have the same shape as they were generated by the IC model itself. This means that the propagation trace of an item *i* must be a sequence of sets of users $D_i(0), \ldots, D_i(n)$, corresponding to the discrete time steps of the IC propagation. Moreover for each node $u \in D_i(t)$ there exists a neighbor *v* of *u* such that $v \in D_i(t-1)$. This is obviously not the case in real-world propagation traces with continuous time.

Following [18] we adopt a delay threshold Δ to define influencers. Specifically, suppose that u adopted i at time $t_i(u)$, and let $t_i(u) = \infty$ if u does not adopt i, then we define $\mathcal{F}_{i,u}^+$ as the set of u' neighbors that potentially influenced u in the selection of i:

$$\mathcal{F}_{i,u}^{+} = \{ v | (v, u) \in E, 0 \le t_i(u) - t_i(v) \le \Delta \}.$$

The set $\mathcal{F}_{i,u}^-$ of u's neighbors who definitely failed in influencing u over i is defined similarly:

$$\mathcal{F}_{i,u}^{-} = \{ v | (v, u) \in E, t_i(u) - t_i(v) > \Delta \}.$$

The main difference between the IC model and TIC, is that while in the former the probability that user v will succeed influencing u is the same for every item i, in the latter $p_{v,u}^i$ is a mixture over the user-to-user influence probabilities, where the mixture weights γ_z^i and the influence probabilities $p_{v,u}^z$ are the parameters to be learned. However, directly unpacking $p_{v,u}$ in order to expose γ_i^z and $p_{v,u}^z$ would lead us to a likelihood formulation which is not tractable in a closed form. We can tackle this problem by resorting to the "complete data" approach [7], which allows us to provide an effective closed form estimation of the parameters γ_i^z and $p_{v,u}^z$. The likelihood of a

propagation trace D_i within the z-th component of the model can be defined as $P(D_i|z;\Theta) = \prod_u P_{u,+}^{i,z} P_{u,-}^{i,z}$, where

$$P_{u,+}^{i,z} = 1 - \prod_{v \in \mathcal{F}_{i,u}^+} (1 - p_{v,u}^z) \quad \text{and}$$
(4)

$$P_{u,-}^{i,z} = \begin{cases} \prod_{v \in \mathcal{F}_{i,u}^-} (1 - p_{v,u}^z) \text{ if } \mathcal{F}_{i,u}^- \neq \emptyset, \\ 1 & otherwise. \end{cases}$$
(5)

We resort to the "complete data" approach [7], by assuming that an unknown binary vector \mathbf{y}_i encodes the information about the generating component. In practice, an item propagates from v to u in a specific (unknown) topic. As a consequence, the complete data likelihood can be defined as

$$\mathcal{L}(\Theta; D_i, \mathbf{Y}) = P(D_i | \mathbf{Y}, \Theta) P(\mathbf{Y} | \Theta)$$
$$= \sum_{z=1}^k y_{iz} P(D_i | z; \Theta) \pi_z$$

In the rest of the paper, following the standard EM notation, $\hat{\Theta}$ will represent the current estimate of the set of parameters Θ . According to the above formulation, the Complete Expectation-Likelihood [7] is given by:

$$\mathcal{Q}(\Theta; \hat{\Theta}) = \sum_{i} \sum_{z=1}^{K} Q_i(z; \hat{\Theta}) \left\{ \log \pi_z + \sum_{u} \log P_{u,+}^{i,z} + \sum_{u} \log P_{u,-}^{i,z} \right\}$$
(6)

Eq. 6 is still untractable in closed form, and it requires a further approximation step in the style of [25]. The trick here is in considering that, if the actual activators are knows, then Eq. 5 can be rewritten in a more tractable format, by assuming a set of bernoulli trials over all possible activators. For each active node v wrt the generic item i, the attemp to activate u succeeds with probability

$$R_{z}^{i}(u,v;\hat{\Theta}) = \frac{\hat{p}_{v,u}^{2}}{\hat{P}_{u,+}^{i,z}}$$
(7)

where \hat{p} and \hat{P} denote the current estimates in $\hat{\Theta}$. $\mathcal{Q}(\Theta; \hat{\Theta})$ can hence be rewritten as

$$\mathcal{Q}(\Theta; \hat{\Theta}) = \sum_{i} \sum_{z=1}^{K} Q_{i}(z; \hat{\Theta}) \left\{ \log \pi_{z} + \sum_{u} \left\{ R_{z}^{i}(u, v; \hat{\Theta}) \log p_{v,u}^{z} + \left(1 - R_{z}^{i}(u, v; \hat{\Theta}) \right) \log(1 - p_{v,u}^{z}) \right\} + \sum_{v \in \mathcal{F}_{i,u}^{-}} \log(1 - p_{v,u}^{z}) \right\}$$

$$(8)$$

Algorithm 1: EM inference of parameters for TIC

```
Input : Social graph G = (V, E), data \mathbb{D}, and K \in \mathbb{N}^+.
Output: The set of all parameters of TIC, \Theta, that is:
                       \forall (v, u) \in E, \forall i \in \mathcal{I}, \forall z \in [1, K] : p_{v, u}^z, \ \pi_z \text{ and } \gamma_i^z.
init(\pi_z, p_{v,u}^z);
repeat
          for
all the i \in \mathcal{I} do
                    for all the z = \{1, \cdots, K\} do
                              Q_{i}(z; \hat{\Theta}) \leftarrow \frac{P(D_{i}|z; \hat{\Theta})\pi_{z}}{\sum_{\bar{z}} P(D_{i}|\bar{z}; \hat{\Theta})\pi_{\bar{z}}};
for all the (u, v) \in E do
      E-step
                                        R_z^i(u,v;\hat{\Theta}) \leftarrow \frac{p_{v,u}^z}{P_{v,u}^{i,z}};
                               end
                    end
          \mathbf{end}
        for all the z = \{1, \cdots, K\} do
                     \begin{aligned} \pi_z &\leftarrow \frac{1}{|\mathcal{I}|} \sum_{i \in \mathcal{I}} Q_i(z; \hat{\Theta}); \\ \text{forall the } (u, v) &\in E : S_{v, u}^+ \neq \emptyset \text{ do} \\ & \Big| p_{v, u}^z \leftarrow \frac{1}{\kappa_{v, u, z}^+ + \kappa_{v, u, z}^-} \sum_{i \in S_{v, u}^+} Q_i(z; \hat{\Theta}) R_z^i(u, v; \hat{\Theta}) \end{aligned} 
       M-step
                    end
          end
until convergence;
```

where π_z is the prior probability that a generic item is assigned to topic z. The mixture parameters γ_i^z , which define the TIC model in Eq. 1, are given by the values of $Q_i(z; \hat{\Theta})$ at the end of the learning procedure. Let $S_{v,u}^+ = \{i | v \in \mathcal{F}_{i,u}^+\}$, and similarly $S_{v,u}^- = \{i | v \in \mathcal{F}_{i,u}^-\}$. Moreover let

$$\kappa_{v,u,z}^+ = \sum_{i \in S_{v,u}^+} Q_i(z; \hat{\Theta}), \text{ and } \kappa_{v,u,z}^- = \sum_{i \in S_{v,u}^-} Q_i(z; \hat{\Theta})$$

The Expectation-Maximization method for learning the parameters of the TIC model is given in Algorithm 1: it starts with a random initialization of parameters π_z (ensuring that $\sum_z \pi_z = 1$) and $p_{v,u}^z$ for all pair $\langle v, u \rangle$ such that $S_{v,u}^+ \neq \emptyset$. Then it alternates the E-step and the M-step, measuring at each iteration the gain of log-likelihood (Eq. 8) w.r.t. the previous iteration. When the gain is below a given threshold, the algorithm has converged.

3.2 Dealing with new items in TIC

TIC model assumes that for each item we are given distribution over the topics and we have shown how to estimate this distribution by log-likelihood maximization. However, an interesting case is to apply the model to a new item never seen before, e.g., when we want to push a new product in the market. In this case we cannot directly apply the parameter estimation procedure described above, since no propagation trace of the new item is available yet. We have to rely on background knowledge about the item. For instance, the marketing expert might directly define the distribution over the topics for the given new item. Alternatively, item features (e.g., genre, price, etc.) might be available, or a small set of initial adopters might have provided tags.

In the most general setting, let us assume that multiple descriptions, in the form of sets of tags from a vocabulary \mathcal{T} , exist for item *i*. Let \mathbf{w}_i denote the bag of tags obtained by joining all the descriptions of *i* and let w_n denote the *n*-th tag in \mathbf{w}_i . Then, we can extend the expected-likelihood formulation in order to take into account tag-assignments for items and maximize their likelihood. Let $\beta_{w_n,k}$ denote the probability of observing the *n*-th tag in the k-mixture: $\beta_{w_n,k} = P(w_n|z_k)$. Assuming that influence probabilities and tags assignments are conditionally independent given the topic, the probability that the trace of item *i* will be generated by the *z*-th component is:

$$P(D_i|z;\Theta) = \prod_{u} P_{u,+}^{i,z} P_{u,-}^{i,z} \prod_{n=1}^{\mathbf{w}_i} \beta_{w_n,k}^{N(w_n,i)}$$

where $N(w_n, i)$ is the number of times that the tag w_n has been assigned to the item *i*. Then, the Complete-Data expectation likelihood becomes:

$$\mathcal{Q}(\Theta; \hat{\Theta})' = \mathcal{Q}(\Theta; \hat{\Theta}) + \sum_{i} \sum_{z=1}^{K} Q(z; i, \hat{\Theta}) \sum_{n=1}^{|\mathbf{w}_i|} N(w_n, i) \log \beta_{w_n, k}$$

and in the M-step we need to update the β distribution as:

$$\beta_{w_n,k} = \frac{1 + \sum_{i \in \mathcal{I}} Q_i(z; \Theta) N(w_n, i)}{|\mathcal{T}| + \sum_{n'=1}^{|\mathcal{T}|} \sum_{i \in \mathcal{I}} Q_i(z; \hat{\Theta}) N(w_{n'}, i)}$$

Since no diffusion trace has been observed yet, it follows that both $\mathcal{F}_{i,u}^+$ and $\mathcal{F}_{i,u}^-$ are empty. In this case, the $P_{u,+}^{i,z}$ and $P_{u,-}^{i,z}$ components equal 1 and the overall probability reduces to the probability of observing the tags within topic z. In addition, γ_i^z can be computed as

$$\gamma_i^z = \frac{\prod_{n=1}^{|\mathbf{w}|} \beta_{w_n, z} \cdot \pi_z}{\sum_{z'} \prod_{n=1}^{|\mathbf{w}|} \beta_{w_n, z'} \cdot \pi_{z'}}.$$
(9)

3.3 Discussion

The traditional IC and LT models suffer various limitations when it comes to apply them in practice. One first limitation is the treatment of time and the consequent need for some discretization, as we have already highlighted in Section 3.1. Another important limitation is the number of parameters. In fact both LT and IC have influence weights (or probabilities) for each pair of connected users. However, having |E| parameters is unsuitable for real-world social networks where the number of edges is usually extremely large (for instance, Facebook nowadays exhibits |E| > 130 billion). The very large number of parameters, on the one hand makes the learning phase computationally prohibitive (the EM-based method needs to update the influence probability associated to each edge in each iteration), and on the other hand it makes the model prone to overfitting.

Goyal *et al.* [10] empirically found that the greedy algorithm with the IC parameters learned with the EM-based method [25], picks as seeds mainly only nodes which perform a very small number of actions, often just one action, and should not be considered as high influential nodes. For instance, in one experiment they found that the first seed selected is a node that in the propagation traces it performs only one action. But this action propagates to 20 of its neighbors. As a result, the EM-based method ends up assigning probability 1 to all the links from that node to all its 20 neighbors, making it a high influence node, so much influential that it results being picked as the first seed by the greedy algorithm.

These limitations are not solved in the topic-aware TIC and TLT models that we have introduced in this section. Indeed, in TIC and TLT we have $K(|E| + |\mathcal{I}|)$ parameters.

In the next section we introduce the AIR (*Authoritativeness-Interest-Relevance*) propagation model, which assumes that social influence depends on a user authority in the context of a given topic and the interest of the user social neighborhood for that topic. This assumption greatly reduces the number of parameters.

4 The "AIR" Propagation Model

The AIR model has the following parameters:

- Authoritativeness of a user in a topic: For each user $v \in V$ and for each topic $z \in [1, K]$, we are given a weight $p_v^z \in \mathbb{R}$ which measures the strength of v's influence on the topic z. A positive value represent *authoritativeness*, i.e., given a topic, the activation of v with respect to an item will influence v's neighbors to select the item as well; on the other hand, negative values model *distrust*, i.e., the activation of v will discourage the activation of her neighbors.
- Interest of a user for a topic: each user u is defined by a distribution ϑ_u over topics: i.e., $\vartheta_u^z = P(Z = z | u)$ denotes the interest of the user u in the topic z and $\sum_{z=1}^{K} \vartheta_u^z = 1$.
- Relevance of an item for a topic: each topic z is defined by a set of weights $\varphi_z \in \mathbb{R}^{|\mathcal{I}|}$, with $\varphi_i^z \in \mathbb{R}$ being the relevance (or selection) weight for the item i in the topic z. Each topic can be hence characterized by the set of the most relevant items. For example, in the topic "Politics" the weight associated with the selection of the "NYT" is expected to be greater than the one corresponding to "Sport Illustrated".

The working principle of AIR is a generalization of the *threshold model* [13]. At the beginning of the process each user u chooses a threshold θ_u uniformly

at random from [0, 1]. At time t, the decision of u to activate for a given item i depends on the influence exerted by her neighbors who have already activated on i (their authoritativeness) and on topic-wise u's interests and i's relevance. In details, at time t user u actives on i iff

$$P(i|u,t) = \sum_{z} P(z|u)P(i|u,z,t) \ge \theta_u$$

where $P(z|u) = \vartheta_u^z$, while P(i|u, z, t) is the following logistic selection function:

$$P(i|u, z, t) = \frac{\exp\left\{\sum_{v \in V} p_v^z f_v(i, u, t) + \varphi_i^z f(i, u, t)\right\}}{1 + \exp\left\{\sum_{v \in V} p_v^z f_v(i, u, t) + \varphi_i^z f(i, u, t)\right\}}$$
(10)

The selection scaling factors $f_v(i, u, t)$ and f(i, u, t) are used to distinguish potential influencers from non influencers $(f_v(i, u, t) = 0 \text{ if } v \notin \mathcal{F}_i(u, t))$ and to potentially relate influence to time. As observed in [12], the likelihood of an item propagating is likely to decay proportionally to time. In particular, it decays at two different levels: locally the influence exerted by v on u for item idecays with the time elapsed from the moment in which v adopted i; globally the interest in the item i decays as i gets older. The adoption of the selection scaling factors in Eq. 10 allows to directly model both local and global temporal decay, e.g. by ensuring that $f_v(i, u, t) \propto (t_i(v) - t)$, and $f(i, u, t) \propto (t_i - t)$.

Compared to the models presented in the previous section, the AIR influence propagation model has only $K(|V| + |\mathcal{I}|)$ parameters. As a result, AIR is a simpler model, more robust to overfitting and still capable of describing influence propagation in an effective way.

4.1 AIR: learning the parameters

The problem of learning the parameters of the AIR model, has the same input of the learning the parameters for TIC, presented in Section 3.1.

Within the generative process, we can assume that for each given item i, a user u picks a topic z by drawing from her own characteristic distribution over the topics ϑ_u (representing her prior interests). Then, for each timestamp t, u activates on i with probability P(i|u, z, t) defined as in Eq. 10. Given the model parameters, we can compute the likelihood of the data as in Eq. 3. Recall that $D_i(t)$ denotes the set of users who selected the item i at time t, while $C_i(t)$ denotes the set of users who selected i by time t. For sake of notation compactness we use the binary indicators $d_i^u(t) = 1$ if $u \in D_i(t)$, and zero otherwise, and $c_i^u(t) = 1$ if $u \in C_i(t)$, and zero otherwise.

Then the Complete-Data Expectation Likelihood is:

$$\mathcal{L}(\Theta; Q) = \sum_{i} \sum_{u} \sum_{z} Q(z; u, i) \{ \log \vartheta_{u}^{z} + \sum_{\underline{t_{i}}}^{\overline{t_{i}}} d_{i}^{u}(t) \log P(i|u, z, t) + (1 - c_{i}^{u}(t)) \log (1 - P(i|u, z, t)) \}$$
(11)

where $\underline{t_i}$ and $\overline{t_i}$ are the timestamp associated respectively with the first and the last observation of the item i.

Each observation $\langle u, i \rangle$ is associated with a state z of the latent variable, modeling the preference of u for i. Also, for the sake of simplicity, we assume that the hidden topic variable is independent from time. This modeling trick simplifies the formulation of the expected likelihood, and provides the following definition for the expected value:

$$Q(z; u, i) = \frac{\hat{\vartheta}_{z}^{u} \prod_{\underline{t}_{i}}^{\overline{t_{i}}} (P(i|u, z, t))^{d_{i}^{u}(t)} \cdot (1 - P(i|u, z, t))^{(1 - c_{i}^{u}(t))}}{\sum_{z'} \hat{\vartheta}_{z'}^{u} \prod_{\underline{t}_{i}}^{\overline{t_{i}}} (P(i|u, z, t))^{d_{i}^{u}(t)} \cdot (1 - P(i|u, z, t))^{(1 - c_{i}^{u}(t))}}$$

Within the EM framework, the ϑ component can be obtained using standard optimization. The remaining parameters are difficult to solve in a closed form, due essentially to the non-linearity of Eq.10. We overcome this limitation by combining the *Improved Iterative Scaling* algorithm [22] and the *Generalized Expectation-Maximization (GEM)* procedure [19].

Essentially, rather than maximizing $\mathcal{L}(\Theta, Q)$, we look for an upgrade Γ of Θ that guarantees

$$\mathcal{L}(\Theta + \Gamma, Q) \ge \mathcal{L}(\Theta, Q)$$

In practice this corresponds to find for each p_v^z an upgrade δ_v^z and for each item *i* an upgrade η_i^z such that the M-step can be defined as $p_v^z \leftarrow p_v^z + \delta_v^z$ and $\varphi_i^z \leftarrow \varphi_i^z + \eta_i^z$.

Notice that the Expectation Log-Likelihood can be expressed as:

$$\mathcal{L}(\Theta, Q) = \sum_{i} \sum_{u} \sum_{z} Q(z; u, i) \vartheta_{z}^{u}$$

$$+ \sum_{i,u,z} Q(z; u, i) \sum_{\underline{t_{i}}}^{\overline{t_{i}}} d_{i}^{u}(t) \left\{ \sum_{v \in \mathcal{V}} p_{v}^{z} f_{v}(i, u, t) + \varphi_{i}^{z} f(i, u, t) \right\}$$

$$- \sum_{i,u,z} Q(z; u, i) \sum_{\underline{t_{i}}}^{\overline{t_{i}}} d_{i}^{u}(t) \log \left\{ 1 + \exp \left\{ \sum_{v \in \mathcal{V}} p_{v}^{z} f_{v}(i, u, t) + \varphi_{i}^{z} f(i, u, t) \right\} \right\}$$

$$- \sum_{i,u,z} Q(z; u, i) \sum_{\underline{t_{i}}}^{\overline{t_{i}}} (1 - c_{i}^{u}(t)) \log \left\{ 1 + \exp \left\{ \sum_{v \in \mathcal{V}} p_{v}^{z} f_{v}(i, u, t) + \varphi_{i}^{z} f(i, u, t) \right\} \right\}$$

$$(12)$$

define

$$a_{i,u,z,t} = \sum_{v \in \mathcal{V}} p_v^z f_v(i, u, t) + \varphi_i^z f(i, u, t)$$

and

$$a'_{i,u,z,t} = \sum_{v \in \mathcal{V}} (p_v^z + \delta_v^z) f_v(i,u,t) + (\varphi_i^z + \eta_i^z) f(i,u,t)$$

Then,

$$\mathcal{L}(\Theta + \Gamma, Q) - \mathcal{L}(\Theta, Q) =$$

$$\sum_{i,u,z} Q(z; u, i) \sum_{\underline{t_i}}^{\overline{t_i}} d_i^u(t) \left\{ \sum_{v \in \mathcal{V}} \delta_v^z f_v(i, u, t) + \eta_i^z f(i, u, t) \right\}$$

$$- \sum_{i,u,z} Q(z; u, i) \sum_{\underline{t_i}}^{\overline{t_i}} d_i^u(t) \log \left\{ \frac{1 + \exp(a'_{i,u,z,t})}{1 + \exp(a_{i,u,z,t})} \right\}$$

$$- \sum_{i,u,z} Q(z; u, i) \sum_{\underline{t_i}}^{\overline{t_i}} (1 - c_i^u(t)) \log \left\{ \frac{1 + \exp(a'_{i,u,z,t})}{1 + \exp(a_{i,u,z,t})} \right\}$$
(13)

Exploiting the inequality $-\log x \geq 1-x$ and ignoring the (positive) constant terms, we obtain:

$$\mathcal{L}(\Theta + \Gamma, Q) - \mathcal{L}(\Theta, Q) \geq \sum_{i,u,z} Q(z; u, i) \sum_{\underline{t_i}}^{\overline{t_i}} d_i^u(t) \left\{ \sum_{v \in \mathcal{V}} \delta_v^z f_v(i, u, t) + \eta_i^z f(i, t) \right\}$$

$$- \sum_{i,u,z} Q(z; u, i) \sum_{\underline{t_i}}^{\overline{t_i}} d_i^u(t) \frac{1 + \exp(a'_{i,u,z,t})}{1 + \exp(a_{i,u,z,t})}$$

$$- \sum_{i,u,z} Q(z; u, i) \sum_{\underline{t_i}}^{\overline{t_i}} (1 - c_i^u(t)) \frac{1 + \exp(a'_{i,u,z,t})}{1 + \exp(a_{i,u,z,t})}$$

$$(14)$$

Notice now that

$$\frac{1 + \exp(a'_{i,u,z,t})}{1 + \exp(a_{i,u,z,t})} = \frac{1}{1 + \exp(a_{i,u,z,t})} + \frac{\exp(a'_{i,u,z,t})}{1 + \exp(a_{i,u,z,t})}$$
$$= (1 - P(i|u, z, t))$$
$$+ P(i|u, z, t) \cdot \exp\left\{\sum_{v \in \mathcal{V}} \delta_v^z f_v(i, u, t) + \eta_i^z f(i, u, t)\right\}$$

which yields

$$\begin{aligned} \mathcal{L}(\Theta + \Gamma, Q) &- \mathcal{L}(\Theta, Q) \geq \\ \sum_{i,u,z} Q(z; u, i) \sum_{\underline{t_i}}^{\overline{t_i}} d_i^u(t) \left\{ \sum_{v \in \mathcal{V}} \delta_v^z f_v(i, u, t) + \eta_i^z f(i, t) \right\} \\ &- \sum_{i,u,z} Q(z; u, i) \sum_{\underline{t_i}}^{\overline{t_i}} d_i^u(t) \left(1 - P(i|u, z, t) \right) \\ &- \sum_{i,u,z} Q(z; u, i) \sum_{\underline{t_i}}^{\overline{t_i}} d_i^u(t) P(i|u, z, t) \cdot \exp\left\{ \sum_{v \in \mathcal{V}} \delta_v^z f_v(i, u, t) + \eta_i^z f(i, t) \right\} \\ &- \sum_{i,u,z} Q(z; u, i) \sum_{\underline{t_i}}^{\overline{t_i}} \left(1 - c_i^u(t) \right) \left(1 - P(i|u, z, t) \right) \\ &- \sum_{i,u,z} Q(z; u, i) \sum_{\underline{t_i}}^{\overline{t_i}} \left(1 - c_i^u(t) \right) P(i|u, z, t) \cdot \exp\left\{ \sum_{v \in \mathcal{V}} \delta_v^z f_v(i, u, t) + \eta_i^z f(i, t) \right\} \end{aligned}$$
(15)

Further, without loss of generality we can assume that the scaling factors are normalized, i.e.,

$$\sum_{v} f_v(u,i) + f(i,u,t) = 1$$

We can hence exploit the Jensen inequality,

$$\exp\left\{\sum_{v\in\mathcal{V}} \delta_v^z f_v(i, u, t) + \eta_i^z f(i, u, t)\right\} \leq \sum_{v\in\mathcal{V}} f_v(i, u, t) \exp\left\{\delta_v^z\right\} + f(i, u, t) \exp\left\{\eta_i^z\right\}$$

which finally allows the lower bound $\mathcal{L}(\Theta + \Gamma, Q) - \mathcal{L}(\Theta, Q) \geq \mathcal{B}(\Gamma, \Theta, Q)$, where

$$\mathcal{B}(\Gamma,\Theta,Q) = \sum_{i,u,z} Q(z;u,i) \sum_{\underline{t_i}}^{\overline{t_i}} d_i^u(t) \left\{ \sum_{v \in \mathcal{V}} \delta_v^z f_v(i,u,t) + \eta_i^z f(i,u,t) \right\} - \sum_{i,u,z} Q(z;u,i) \sum_{\underline{t_i}}^{\overline{t_i}} d_i^u(t) \left(1 - P(i|u,z,t)\right) - \sum_{i,u,z} Q(z;u,i) \sum_{\underline{t_i}}^{\overline{t_i}} d_i^u(t) P(i|u,z,t) \cdot \cdot \left\{ \sum_{v \in \mathcal{V}} f_v(i,u,t) \exp\left\{\delta_v^z\right\} + f(i,u,t) \exp\left\{\eta_i^z\right\} \right\}$$
(16)
$$- \sum_{i,u,z} Q(z;u,i) \sum_{\underline{t_i}}^{\overline{t_i}} \left(1 - c_i^u(t)\right) \left(1 - P(i|u,z,t)\right) - \sum_{i,u,z} Q(z;u,i) \sum_{\underline{t_i}}^{\overline{t_i}} \left(1 - c_i^u(t)\right) P(i|u,z,t) \cdot \cdot \left\{ \sum_{v \in \mathcal{V}} f_v(i,u,t) \exp\left\{\delta_v^z\right\} + f(i,u,t) \exp\left\{\eta_i^z\right\} \right\}$$

The parameter update Γ improves the likelihood when the auxiliary function $\mathcal{B}(\Gamma, \Theta, Q)$ is positive. By maximizing the latter, we obtain:

$$\begin{split} \delta_{v}^{z} &= \log \left\{ \frac{\sum_{i,u} Q(z;u,i) f_{v}(i,u,t_{i}(u))}{\sum_{i,u} Q(z;u,i) \sum_{\underline{t_{i}}}^{t_{i}(u)} P(i|u,z,t) \cdot f_{v}(i,u,t)} \right\} \\ \eta_{i}^{z} &= \log \left\{ \frac{\sum_{u} Q(z;u,i) f(i,u,t_{i}(u))}{\sum_{u} Q(z;u,i) \sum_{\underline{t_{i}}}^{t_{i}(u)} P(i|u,z,t) \cdot f(i,u,t)} \right\} \end{split}$$

Algorithm 2 summarizes the overall learning scheme.

4.2 AIR: dealing with new items

.

Modeling unobserved items follows the general guidelines exposed in Sec. 3.2, with some variations. The selection probability for a new item can be simplified as:

$$P(i|u,t) = \sum_{z} P(z|u)P(i|z,u,t)$$

=
$$\sum_{z} \vartheta_{u}^{z} \frac{\exp\left\{\varphi_{i}^{z}f(i,u,t)\right\}}{1 + \exp\left\{\varphi_{i}^{z}f(i,u,t)\right\}}$$
(17)

Algorithm 2: EM inference of parameters for AIR **Input** : Social graph G = (V, E), data \mathbb{D} , and $K \in \mathbb{N}^+$. $\mathbf{Output}:$ The set of all parameters of $\mathsf{AIR}\ \varTheta,$ that are $p_u^z(\mathsf{A}), \vartheta_u^z(\mathsf{I}), \varphi_i^z(\mathsf{R}), \text{forall } u \in V, z \in [1, K], i \in \mathcal{I}.$ $init(p_u^z, \vartheta_u^z, \varphi_i^z); //Random initialization of parameters$ repeat for all the $i \in \mathcal{I}$ do for all the $u \in V$ do forall the $z = \{1, \dots, K\}$ do E-step $Q(z;u,i) \leftarrow \frac{\vartheta_{u}^{z} \prod_{t_{i}}^{\overline{t_{i}^{z}}} P(i|u,z,t)^{d_{i}^{u}(t)} \cdot (1 - P(i|u,z,t))^{(1 - c_{i}^{u}(t))}}{\sum_{z'} \vartheta_{u}^{z'} \prod_{t_{i}}^{\overline{t_{i}^{z}}} P(i|u,z',t)^{d_{i}^{u}(t)} \cdot (1 - P(i|u,z',t))^{(1 - c_{i}^{u}(t))}}$ end end end forall the $z = \{1, \cdots, K\}$ do forall the $v \in V$ do $\vartheta_v^z \leftarrow \frac{1}{|\mathcal{I}|} \sum_{i \in \mathcal{I}} Q(z; v, i)$ $\delta_v^z \leftarrow \log\left\{\frac{\sum_{i,u} Q(z;u,i)f_v(i,u,t_i(u))}{\sum_{i,u} Q(z;u,i)\sum_{t_m(i)}^{t_i(u)} P(i|u,z,t) \cdot f_v(i,u,t)}\right\}$ \mathbf{end} for all the $i \in \mathcal{I}$ do $\eta_i^z \leftarrow \log \left\{ \frac{\sum_u Q(z;u,i)f(i,t_i(u))}{\sum_u Q(z;u,i)\sum_{t_m(i)}^{t_i(u)} P(i|u,z,t) \cdot f(i,t)} \right\}$ M-step end for all the $u \in V$ do $| p_u^z \leftarrow p_u^z + \delta_u^z$ end for all the $i \in \mathcal{I}$ do $\varphi_i^z \leftarrow \varphi_i^z + \eta_i^z$ end end until convergence;

In general relevance parameter φ_i^z for a new item *i* is not bound to an optimal value. However, when tag information is available, we can assume a prior tendency of the item to be selected, according to its likelihood to be associated with the topic. That is, we can model new item by assuming a prior probability $p(\varphi_i^z)$, defined as a gaussian distribution with constant variance σ and mean γ_i^z (as defined in Eq. 9). As a consequence, the log-likelihood can be reformulated to comprise the prior probabilities, resulting into

$$\mathcal{L}(\Theta; \mathbb{D}) = \sum_{i \in \mathcal{I}} \log \mathcal{L}(\Theta; D_i) + \log P(\Theta)$$

A more thorough maximum a posteriori estimation (MAP) treatment for the whole parameter set Θ is omitted here for lack of space. Without loss of generality, we assume uniform prior probabilities for all the parameters other than new items. As a consequence, $P(\Theta)$ can be simplified as:

$$\log P(\Theta) = \sum_{i:D_i = \emptyset} \sum_{z} \log p(\varphi_i^z) + C$$

Combining the above equations and Eq. 11 finally yields

$$\mathcal{L}(\Theta; Q)' = \mathcal{L}(\Theta; Q) + \sum_{i: D_i = \emptyset} \sum_{z} \log p(\varphi_i^z) + C$$

Optimizing the latter with respect to a parameter φ_i^z relative to a new item *i* yields the straightforward solution $\varphi_i^z = \gamma_i^z$.

4.3 Influence Maximization in AIR

We next discuss the problem of influence maximization in AIR. Given a generic item *i* that we want to promote, we assume that its AIR parameters are known. The problem is to select a set *S* of *k* nodes such that the expected spread of influence of *S* under the AIR model, denoted $\sigma_{AIR}(S)$, is maximal.

Although AIR is a general threshold model, the fact that user authoritativeness can be negative makes σ_{AIR} not submodular and not even monotone. Therefore the standard greedy algorithm cannot provide any approximation guarantee, as it does for the classic IC and LT models, and for their topic-aware versions TIC and TLT.

Even without any provable guarantee, it is reasonable to consider the greedy algorithm a reasonable candidate also for the AIR model, given that in any case we shall naturally avoid users with negative authoritativeness. Therefore, in the next section we compare the spread $\sigma_{AIR}(S)$ achieved by the following two methods:

- **Greedy:** at each iteration greedily add to the set of seeds S the node x that brings the largest marginal gain, i.e., $\sigma_{AIR}(S \cup \{x\}) \sigma_{AIR}(S)$ is maximal. Estimate $\sigma_{AIR}(S)$ for a given S by Monte Carlo simulations [13].
- **Top-**k **authorities:** given the new item i and its distribution over topics γ_i^z , select the top-k users v w.r.t.

$$\sum_{z=1}^{K} \gamma_i^z p_v^z.$$

Recall that all over the paper K is the number of topics, while here k is the size of the required seed set.

Studying alternative approaches to influence maximization under the AIR model will be part of our future investigation.

	FLIXSTER		Digg		
	Training	Test	Training	Test	
Users	6,572	4,686	16,297	14,061	
Items	7,158	7,138	3,553	3,547	
Actions	1,432,716	340,495	1,160,428	264,066	
Avg $\#$ actions (user)	218	72	71	18	
Avg $\#$ actions (item)	200	47	326	74	
Min $\#$ actions (user)	6	1	6	1	
Min $\#$ actions (item)	9	1	90	3	
Max # actions (user)	5,525	1,786	$2,\!640$	1,912	
Max # actions (item)	3,173	778	4,995	828	
Avg lifetime (item)	952 days		14 days		
Avg time between two actions					
per user	94 hours		66 hours		
per item	22 days		38 minutes		

Table 2: Summary of the propagation data.

5 Experimental Evaluation

The goal of our experiments is twofold. At a high level, we want to evaluate the impact of introducing a topic-based estimation of the influence probabilities. That is, we are interested in evaluating whether topic-aware propagation models can better predict the activation of a user on a specific item. The expected result is that the combined adoption of both influence and topic modeling exhibits an improvement over the single contributions. We also aim aim at assessing whether considering the topic model of the item can bring any benefit in a viral marketing campaign. That is to say, to compare topic-aware models against models that ignore the topic distribution of the item, in the influence maximization problem.

5.1 Datasets

We use two real-world and publicly available datasets, both containing a social graph G = (V, E) and a log of past propagations $\mathbb{D} = \{ (User, Item, Time) \}$: the datasets come from Digg (www.digg.com) and Flixster (www.flixster.com). Digg is a social news website, where the users vote stories. In this case \mathbb{D} contains information about which user voted which story (item) at which time. If we have user v vote a story about the new iPhone, and shortly later v's friend u does the same, we consider the story as having propagated from v to u, and v as a potential influencer for u. Flixster is one of the main players in the mobile and social movie rating business. Here, an item is a movie, and the action of the user is rating the movie.

In both cases we started from the publicly available dataset 2 3 and we performed some standard consistency cleaning and removal of all users and

² www.isi.edu/~lerman/downloads/digg2009.html

³ http://www.cs.sfu.ca/~sja25/personal/datasets/



Fig. 1: Frequency distributions for influencers

items that do not appear at least 20 times in \mathbb{D} . The final DIGG social graph contains 11,142 users and 99,846 directed arcs, while FLIXSTER contains 6,353 users and 84,606 directed arcs: in both cases we do not consider the disconnected nodes, i.e., users that appear actively in \mathbb{D} but which have no friends in *G*. Moreover, for our purposes we performed a chronological split of \mathbb{D} in both datasets into training (80%) and test (20%). Table 2 summarizes the main properties of \mathbb{D} .

The two datasets exhibit different features and it is worth analyzing the distribution of the number of potential influencers for each activation in their respective propagation log \mathbb{D} . More specifically, for each tuple (u, i, t) we record $|\mathcal{F}_{i,u}^+|$; the cumulative distributions of the overall number of potential influencers are given in Figure 1. This analysis allows us to measure the role, and extent, of influence in the behavior of the users on the two datasets. We can observe that influence plays a more relevant role in the Flixster dataset, and we expect that the next analysis of the influence weights will confirm this hypothesis.

5.2 Experiments settings

We start by noticing that there is a direct relationship between the the scaling factors $f_v(i, u, t)$ of the AIR model and the size of influence window Δ used in the parameters learning of the IC and TIC models. We studied two alternative definitions for $f_v(i, u, t)$.

The first one assumes that propagation can degrade following an exponential decay:

$$f_{v}(i, u, t) \propto \begin{cases} \exp(t_{v}(i) - t) \text{ if } v \in \mathcal{F}_{i}(u, t) \\ 0 & \text{otherwise} \end{cases}$$



Fig. 2: Convergence rate: AIR vs. TIC on DIGG (left) and FLIXSTER (right).

This definition of scaling factor corresponds to a very short influence threshold Δ (typically, 3 to 5 timestamps). The second option we explored is to keep $f_v(i, u, t)$ constant. This corresponds to adopting a value $\Delta = \infty$ within the IC and TIC models. As a matter of fact, the statistics on the average time between two actions involving the same item, and the average time-life for an item in Table 2 suggest for a large Δ . Our empirical analysis determined that, at least in these two datasets, the best results are achieved by considering all the influencers up to the considered time: i.e., $\Delta = \infty$ and consequently $f_v(i, u, t)$ constant.⁴ Therefore in the experiments reported here we always adopt these settings.

5.3 Learning Rate

In Figure 2 we compare the learning rate of the TIC and the AIR model in the first 200 iterations. As expected, TIC exhibit a faster convergence rate than AIR: this is due to the difference in their respective M-step. AIR relies on a *GEM* procedure which clearly affects the number of iterations needed to achieve convergence. Notably, the TIC parameter estimation phase provides a good estimation of the model parameters after about 60 iterations on both the datasets, whereas the AIR model requires approximately 1400. Also, both algorithms are initialized randomly, but the likelihood increase for AIR is slower. We plan to investigate ways to speed up the parameter estimation phase of the AIR model, as well as better initialization strategies in future works.

 $^{^4~}$ This is in accordance with the experiments in [18], that firstly introduced the \varDelta influence window.

5.4 Analysis of the influence weights

The analysis of the influence weights that characterize the proposed topicaware propagation models can provide meaninful insights about the users' behavior in the considered datasets. More specifically, we are going to focus on the following aspects:

- How are the influence weights distributed?
- Do the main differences in the two considered datasets reflect on the distributions of the influence weights?
- Does users' authoritativeness/influence change on different topics?

Figure 3 plots the distributions of the influence weights of the AIR model for the number of topics achieving the best performances on the two considered datasets (as described later in this section). Values are distributed according to two log-normal distributions centered in the positive and negative quadrants, with relatively slow values and relatively few extreme values. The graphs show that negative influencers also play a significant role in the learning phase. The bottom graphs also show how the item-topic weights distribute. Again, it seems that the AIR model has a bias towards negative weights. The Digg dataset exhibits a lower level of influence among users, as witnessed by the high number of weights set to 0. This difference between the two datasets can be explained by considering the distribution of the number of influencers for each activation, provided in Figure 1. Figure 4 plots of the influence probabilities for the TIC models, and confirm such a trend. Here, values are exponentially distributed; however a percentage of users exhibit highest influence.

A key point that motivated the introduction of the AIR model was that users' may exert different degree on authoritativeness on different topics. To verify the correctness of such assumption, we can analyze the variance of the user's influence weights on different topics. Figure 5 plots the cumulative frequencies of users' standard deviations on influence weights, both for TIC and AIR. In the plots, only users with positive deviation were considered. For both models, we can observe an higher variance on Flixster, which confirms the intuition that Flixster is more susceptible to influence than Digg.

Figure 6 shows the distributions of frequencies of positive weights and negative weights for a given user. Again, there is a tendency to exhibit more negative values. This is also witnessed by the graph that shows the difference between the number of positive and negative in a user (a negative value here denotes that there are more negative weights). Notwithstanding, some users tend to exhibit a predominance of positive weights. This is a clue that the model learns to discriminate between positive and negative influences. Also, it is worth noticing that in the majority of cases there is a mix of positive and negative weights, clearly stating that the influence of various users changes depending on the topics.

We also analyze, in Figure 7, the correlation between the users's authoritativeness score and the out degree, by distinguishing positive and negative average authoritativeness values. In both the cases, we do not register any sig-



Fig. 3: Distribution of $p_v^z(\text{first row})$ and φ_i^z (secondo row) in the AIR model;



Fig. 4: Distribution of $p_{v,u}^{z}$ in the TIC model (non-zero values).

nificant correlation value that could support the hypothesis for which high/low connected users are also high/low influential.



Fig. 5: Variance of the influence weights: (a) TIC; (b) AIR.



Fig. 6: Distribution of frequencies of positive/negative weights, and difference between them. Digg (first row) and Flixster (second row).

Finally, Figure 8 shows how the interest of a user for a topic varies. Deviation is slow in Digg (an effect of the high number of topics as well). However, it is interesting to notice that Flixster exhibit an high average deviation, and in general the low values are rare.



Fig. 7: Authoritativeness vs Out Degree in AIR models on Digg (first row), and Flixster (second row).



Fig. 8: Variance of the users' interests in the AIR models.



Fig. 9: ROC analysis: DIGG (left column) and Flixster (right column).

5.5 Predictive accuracy

In the following we compare IC, TIC and AIR: the parameters of the model are learned using the EM method in [25], the method in Section 3.1, and in Section 4.1 respectively. The basic principle guiding our evaluation can be summarized as follows. Given the training propagation data \mathbb{D}_T and a test propagation data \mathbb{D}_{Test} , a generic model, whose parameters have been learned on \mathbb{D}_T , provides a suitable estimation of influence and behavior if its application to unobserved data \mathbb{D}_{Test} provides accurate predictions, which can be measured through the following tests.

Activation Test (General). The idea is to measure whether a diffusion model can predict the overall user's activations. This is basically a binary prediction task: for a given user-item pair $\langle u, i \rangle \notin \mathbb{D}_T$, we try to predict whether $\langle u, i \rangle \in \mathbb{D}_{Test}$. Since this test is time-independent, we also use as a baseline for comparison the *Probabilistic Latent Semantic Analysis (pLSA)* model [11]. Although not originally aimed at modeling influence, the latter also relies on topic modeling and occurrences of user actions. Hence, its inclusion in the test allows us to evaluate the contribution of topic modeling on the activation prediction.

Selection Probabilities (General). For each pair $\langle u, i \rangle$ we measure the degree of responsiveness of the model at the actual activation time $t_i(u)$ (if it exists). A good model should assign high probability of activation to a positive case $\langle u, i \rangle \in \mathbb{D}_{Test}$, and low probability (relative to all the possible timestamps) to a negative case $\langle u, i \rangle \notin \mathbb{D}_{Test}$, while the model should assigns low probability for each considered timestamp if the user has not selected the considered item. More specifically, for each pair $\langle u, i \rangle \notin \mathbb{D}_T$ we compute the probability $P(i|u, t'_i(u))$, where

$$t'_{i}(u) = \begin{cases} t_{i}(u) & \text{if } \langle u, i \rangle \in \mathbb{D}_{Test} \\ \arg \max_{t \in [t_{i}, \overline{t_{i}}]} P(i|u, t) & \text{otherwise} \end{cases}$$

Selection Probabilities (Influence Episodes). The previous test strongly penalizes pure influence-based diffusion models, as they assign zero probability to episodes $\langle u, i, t \rangle$ for which the set of influencers is empty. This is not true, of course, for the AIR model which is able to capture the relevance of an item for a given topic. In order to measure the effects of influence, in this test we focus on those episodes $\langle u, i, t \rangle$ for which $\mathcal{F}_i(u, t) \neq \emptyset$.

Activation Time (Influence Episodes). A final test measures the precision of activation at a given timestamp, only considering episodes with non-empty influencess set. Each pair $\langle u, i \rangle \notin \mathbb{D}_T$ is evaluated by comparing the true activation time (if any) with the predicted activation time. Let $t'_i(u)$ represent the predicted activation timestamp, i.e., the minimal timestamp t where P(i|u, t)is greater than a given activation threshold (with $t'_i(u) = \infty$ if the model does not indeed predict any activation for the given item). We can devise the following confusion matrix:

	$\langle u, i \rangle \in \mathbb{D}_{Test}$	$\langle u, i \rangle \not\in \mathbb{D}_{Test}$
True Positive	$t_i'(u) = t_i(u)$	-
False Positive	$t_i'(u) < t_i(u)$	$t'_i(u) \neq \infty$
True Negative	-	$t'_i(u) = \infty$
False Negative	$t_i'(u) > t_i(u)$	-

For all the above mentioned tests we plot the *Receiver Operating Characteristic* (ROC) curves relative to varying activation thresholds. The results are given in Figure 9, while in Table 3 we report the *Area Under the Curve* (AUC) values.

Evaluation. We experimentally found that the optimal number of topics on Digg is 15 topics for AIR, and 20 on TIC. Also, Flixster settles 3 topics on AIR, and 10 on TIC.

The AIR models achieve the best results in detecting the activations, with a consistent gain over the other models (including the runner-up pLSA model). Independent cascade models (IC and TIC) exhibit partial curves on this test, limiting the upper bound of FPR to 0.1. This is due to the fact that negative cases $\langle u, i \rangle$ are a vast majority, and when the case exhibit no active influencers the IC and TIC models assign 0 probability, which eventually results in a True Negative in the test. For this reason, the extension to topics does not provide a significant improvement: both IC and TIC overlap, and the difference in AUC is marginal. Things change when activation time is taken into account: in the remaining plots, TIC outperforms IC, an evidence of a substantial contribution of the topic modeling in increasing the accuracy of time-oriented predictions. Again, AIR achieves the best accuracy among all the models under investigation.

Tests in Fig.9(d) are the most fine-grained: here underestimation of influence (resulting in retarded activation prediction) as well as overestimation (resulting in anticipated activation prediction) are paid as errors. Clearly, topic modeling plays a crucial role in this test, as it allows to better correlate the estimation phase to the actual activation time.

5.6 Influence Maximization

We now turn our attention to the influence maximization problem and to the following questions: (1) how important is it to consider the topic-distribution of the item while selecting the seed sets? (2) how good are the greedy algorithm and the top-k-authorities heuristic on the AIR model? (3) how much does the item "popularity" affect the overall spread?

In Figure 10(left) we compare the expected spread achieved on the AIR propagation model by the greedy algorithm and the top-k-authorities heuristic. The experiment is performed on FLIXSTER, using 50 different items, and averaging the results. Items are described by their relevance over 3 topics. We also add to the comparison a seed set selected by the greedy algorithm on the IC model: i.e., without considering the topics. Being topic-blind, the IC experiment is run only for one generic item. All the greedy algorithms use 1000 Monte Carlo simulations to estimate the expected spread.

	Digg	Flixster	Digg	FLIXSTER	
Model	Activation Test (General)		Selection Probs. (General)		
AIR	0.8585511	0.8857634	0.8484368	0.8201586	
TIC	0.6190136	0.731208	0.6256339	0.7000218	
IC	0.6189209	0.730694	0.5256555	0.702175	
	Selection Probs. (Inf. episodes)		Activation Time		
AIR	0.8123432	0.783486 4	0.8784483	0.8150082	
TIC	0.7714797	0.7253222	0.7377654	0.7377654	
IC	0.7916101	0.6940882	0.6294611	0.6089646	

Table 3: Summary of the evaluation: AUC values.



Fig. 10: Influence maximization experiments.

Although (as discussed in Section 4.3) the greedy algorithm does not provide approximation guarantee, it outperforms the top-k-authorities heuristic. The latter still performs very well: it achieves a spread quite close to those of the greedy approach, and in addition it is much faster to compute. More importantly, both topic-aware strategies largely outperform the topic-blind IC-greedy strategy.

In Figure 10 (right) we compare a "popular" item, i.e., an item which has a rather high relevance (a value of 10) in all three topics, with a normal item having relevance 10 in one topic, and relevance 1 in the other two topics. Not surprisingly, we can observe that the popular item achieves a larger spread. The difference tends to decrease with larger seedset: apparently, popular items are tolerant to smaller seedsets, whereas general items require more seeds.

6 Conclusions and Future work

We provided a topic-modeling perspective over social influence, by introducing novel topic-aware propagation models. We devised methods to learn model parameters from a log of past propagations. We experimentally found the proposed models more accurate in describing real-world influence-driven propagations than the state-of-the-art approaches: as a matter of fact, the two proposed models exhibit an average 28% (AIR) and 7% (TIC) improvement on AUC over the baseline IC approach. The tests show that the models provide accurate predictions of both activations and activation times, and they provide robust estimates of influence parameters. The latter can be used to explain the factors influencing a user's action in a social network, and ultimately to predict their behavior and preferences. Finally, we showed that by considering the characteristics of the item we can obtain larger spread in influence maximization.

There are several ways to extend the main results of this paper. First of all, we plan to investigate ways to speed up the parameter estimation phase of the AIR model, as well better initialization. Also, from a modeling perspective, a full bayesian treatment of the topic models introduced here can help with model generalization and overfitting avoidance.

We also plan to study influence maximization methods based on the AIR model. Finally, we plan to extend the focus of this paper to further application domains, by investigating how to combine influence maximization with topic modeling for recommender systems.

Acknowledgments. This research was partially supported by the Torres Quevedo Program of the Spanish Ministry of Science and Innovation, and partially funded by the European Union 7th Framework Programme (FP7/2007-2013) under grant n. 270239 (ARCOMEM).

References

- Bakshy, E., Hofman, J.M., Mason, W.A., Watts, D.J.: Everyone's an influence: quantifying influence on twitter. In: Proc. of the Forth Int. Conf. on Web Search and Web Data Mining (WSDM'11) (2011)
- Blei, D.M.: Introduction to probabilistic topic models. Communications of the ACM (2011). URL http://www.cs.princeton.edu/~blei/papers/Blei2011.pdf
- Blei, D.M., Ng, A.Y., Jordan, M.I.: Latent dirichlet allocation. The Journal of Machine Learning Research 3, 993–1022 (2003)
- 4. Bonchi, F.: Influence propagation in social networks: A data mining perspective. IEEE Intelligent Informatics Bulletin, Vol.12 No.1: 8-16 (2011)
- Chen, W., Wang, C., Wang, Y.: Scalable influence maximization for prevalent viral marketing in large-scale social networks. In: Proc. of the 16th ACM SIGKDD Int. Conf. on Knowledge Discovery and Data Mining (KDD'10) (2010)
- Chen, Y.C., Peng, W.C., Lee, S.Y.: Efficient algorithms for influence maximization in social networks. Knowl. Inf. Syst. 33(3), 577–601 (2012)
- Dempster, A.P., Laird, N.M., Rubin, D.B.: Maximum Likelihood from Incomplete Data via the EM Algorithm. Journal of the Royal Statistical Society. Series B (Methodological) 39, 1–38 (1977)
- Domingos, P., Richardson, M.: Mining the network value of customers. In: Proc. of the Seventh ACM SIGKDD Int. Conf. on Knowledge Discovery and Data Mining (KDD'01) (2001)
- Goyal, A., Bonchi, F., Lakshmanan, L.V.S.: Learning influence probabilities in social networks. In: Third ACM Int. Conf. on Web Search and Data Mining (WSDM'10) (2010)
- Goyal, A., Bonchi, F., Lakshmanan, L.V.S.: A data-based approach to social influence maximization. PVLDB 5(1), 73–84 (2011)

- Hofmann, T.: Probabilistic Latent Semantic Analysis. In: Proceedings of Uncertainty in Artificial Intelligence, UAI (1999)
- Ienco, D., Bonchi, F., Castillo, C.: The meme ranking problem: Maximizing microblogging virality. In: Proc. of the SIASP workshop at ICDM'10 (2010)
- Kempe, D., Kleinberg, J.M., Tardos, É.: Maximizing the spread of influence through a social network. In: Proc. of the Ninth ACM SIGKDD Int. Conf. on Knowledge Discovery and Data Mining (KDD'03) (2003)
- Kimura, M., Saito, K.: Tractable models for information diffusion in social networks. In: J. Frnkranz, T. Scheffer, M. Spiliopoulou (eds.) Knowledge Discovery in Databases: PKDD 2006, Lecture Notes in Computer Science, vol. 4213, pp. 259–271. Springer Berlin / Heidelberg (2006). URL http://dx.doi.org/10.1007/11871637_27
- Leskovec, J., Krause, A., Guestrin, C., Faloutsos, C., VanBriesen, J., Glance, N.S.: Costeffective outbreak detection in networks. In: Proc. of the 13th ACM SIGKDD Int. Conf. on Knowledge Discovery and Data Mining (KDD'07) (2007)
- Lin, X., Mei, Q., Han, J., Jiang, Y., Danilevsky, M.: The joint inference of topic diffusion and evolution in social communities. In: Proc. of IEEE Int. Conference on Data Mining (ICDM'11) (2011)
- Liu, L., Tang, J., Han, J., Jiang, M., Yang, S.: Mining topic-level influence in heterogeneous networks. In: Proc. of the 19th ACM Conf. on Information and Knowledge Management (CIKM'10) (2010)
- Mathioudakis, M., Bonchi, F., Castillo, C., Gionis, A., Ukkonen, A.: Sparsification of influence networks. In: Proc. of the 17th ACM SIGKDD Int. Conf. on Knowledge Discovery and Data Mining (KDD'11) (2011)
- Neal, R.M., Hinton, G.E.: A view of the EM algorithm that justifies incremental, sparse, and other variants, pp. 355–368. MIT Press, Cambridge, MA, USA (1999)
- Nemhauser, G.L., Wolsey, L.A., Fisher, M.L.: An analysis of approximations for maximizing submodular set functions i. Mathematical Programming 14(1), 265–294 (1978)
- Nicola Barbieri Francesco Bonchi, G.M.: Topic-aware social influence propagation models. In: Proc. of IEEE Int. Conference on Data Mining (ICDM'12) (2012)
- Nigam, K.: Using maximum entropy for text classification. In: In IJCAI-99 Workshop on Machine Learning for Information Filtering, pp. 61–67 (1999)
- Richardson, M., Domingos, P.: Mining knowledge-sharing sites for viral marketing. In: Proc. of the Eighth ACM SIGKDD Int. Conf. on Knowledge Discovery and Data Mining (KDD'02) (2002)
- Saito, K., Kimura, M., Ohara, K., Motoda, H.: Efficient discovery of influential nodes for sis models in social networks. Knowl. Inf. Syst. 30(3), 613–635 (2012)
- Saito, K., Nakano, R., Kimura, M.: Prediction of information diffusion probabilities for independent cascade model. In: Proc. of the 12th Int. Conf. on Knowledge-Based Intelligent Information and Engineering Systems (KES'08) (2008)
- Song, X., Chi, Y., Hino, K., Tseng, B.L.: Information flow modeling based on diffusion rate for prediction and ranking. In: Proc. of the 16th Int. Conf. on World Wide Web (WWW'07) (2007)
- Steyvers, M., Griffiths, T.: Probabilistic topic models. In: T. Landauer, D. Mcnamara, S. Dennis, W. Kintsch (eds.) Latent Semantic Analysis: A Road to Meaning. Laurence Erlbaum (2006). URL http://cocosci.berkeley.edu/tom/papers/ SteyversGriffiths.pdf
- Tang, J., Sun, J., Wang, C., Yang, Z.: Social influence analysis in large-scale networks. In: Proc. of the 15th ACM SIGKDD Int. Conf. on Knowledge Discovery and Data Mining (KDD'09) (2009)
- Weng, J., Lim, E.P., Jiang, J., He, Q.: Twitterrank: finding topic-sensitive influential twitterers. In: Proc. of the Third Int. Conf. on Web Search and Web Data Mining (WSDM'10) (2010)
- Xiang, R., Neville, J., Rogati, M.: Modeling relationship strength in online social networks. In: Proc. of the 19th Int. Conf. on World Wide Web (WWW'10) (2010)

Author Biographies

nicola.pdf

francesco.pdf

beppe.pdf

Nicola Barbieri is a post-doc in the the Web Mining Research Group at Yahoo! Research Barcelona. He graduated with full marks and honor and received his PhD in 2012 at University of Calabria - Italy, and he was a fellow researcher at ICAR-CNR. His research focuses mainly on probabilistic topic models with applications on recommender systems and collaborative filtering, and social influence information propagation in social networks.

Francesco Bonchi is a senior research scientist at Yahoo! Research in Barcelona, Spain, where he is the head of the Web Mining Research Group. His recent research interests include mining query-logs, social networks, and social media, as well as the privacy issues related to mining these kinds of sensible data. He is member of the ECML PKDD Steering Committee and he serves in the editorial board of IEEE Transactions on Knowledge and Data Engineering (TKDE) and ACM Transactions on Intelligent Systems and Technology (TIST). More information can be found at http://www.francescobonchi.com/.

puter science from the University of Pisa. He is currently a senior researcher at the Institute of High Performance Computing and Networks (ICAR-CNR) of the National Research Council of Italy and a contract professor at University of Calabria, Italy. He has been contract researcher at the CNUCE Institute in Pisa, Italy, and a visiting fellow at the CWI Institute in Amsterdam, Netherlands. His current research interests include knowledge discovery and data mining, Recommender systems and Social Network analysis.

Giuseppe Manco Giuseppe Manco graduated summa cum laude in computer science and received the PhD degree in com-

