

# The Joint Inference of Topic Diffusion and Evolution in Social Communities

Cindy Xide Lin<sup>1</sup>   Qiaozhu Mei<sup>2</sup>   Jiawei Han<sup>1</sup>   Yunliang Jiang<sup>1</sup>   Marina Danilevsky<sup>1</sup>

<sup>1</sup>Department of Computer Science, University of Illinois at Urbana-Champaign,

<sup>2</sup>School of Information, University of Michigan,

xidelin2@uiuc.edu,   qmei@umich.edu,   hanj@cs.uiuc.edu,   jiang8@uiuc.edu,   danilev1@illinois.edu

**Abstract**—The prevalence of Web 2.0 techniques has led to the boom of various online communities, where topics spread ubiquitously among user-generated documents. Working together with this diffusion process is the evolution of topic content, where novel contents are introduced by documents which adopt the topic. Unlike explicit user behavior (e.g., buying a DVD), both the diffusion paths and the evolutionary process of a topic are implicit, making their discovery challenging.

In this paper, we track the evolution of an arbitrary topic and reveal the latent diffusion paths of that topic in a social community. A novel and principled probabilistic model is proposed which casts our task as an joint inference problem, which considers textual documents, social influences, and topic evolution in a unified way. Specifically, a mixture model is introduced to model the generation of text according to the diffusion and the evolution of the topic, while the whole diffusion process is regularized with user-level social influences through a Gaussian Markov Random Field. Experiments on both synthetic data and real world data show that the discovery of topic diffusion and evolution benefits from this joint inference; and the probabilistic model we propose performs significantly better than existing methods.

## I. INTRODUCTION

The prevalence of Web 2.0 techniques has led to the boom of various online communities. One of the core problems with analyzing such online communities is concerned with understanding the cascading behaviors and the diffusion of information. Epidemic diseases, adoption of innovation, information memes, and many types of user actions all spread widely in these communities, following the social network of users. The modeling of information diffusion plays a crucial role in many domains. The contagion of disease forms the foundation of epidemics; the social influence in cascading behaviors has been a basic mechanism of viral marketing; and the diffusion of topics is essential to the understanding of scientific innovation. Recently, a large body of research work has been done in the field of social network analysis, aiming to describe the macro-level dynamics and characteristics of information diffusion [17], [11], revealing key factors that affect the adoption of behaviors [1], [18], and describing the design of contagion models that simulate the diffusion process [25], [35].

Many conclusions in this line of work are motivated and validated in scenarios where the actual contagion/diffusion paths are observed. However, such an assumption, which is considered to be common practice in user surveys and

controlled user studies, does not apply to large scale online communities. While the adoptions of behaviors are relatively easy to observe (based on which most macro-level descriptive statistics are computed), the evidence of actual contagion and influence tend to be vague. Who infected whom? Who got the gossip from whom? Who influenced whose research? There are still substantial challenges in this micro-level analysis of information diffusion in large scale social networks. Indeed, users who joined a community or purchased an iPad usually won't explain which particular friends have influenced them; rumor spreaders tend to cover the source of the information; and a researcher cites many references in her paper, without, for instance, labeling the top three who have had the most salient influence on her work. The identification of contagion is difficult even if the general social network structure is observed. It is a non-trivial task to detect the actual diffusion paths of user behaviors merely based on the time of adoption and the social network structure, known as the problem of diffusion inference (or influence) [8].

Inference of diffusion [1], [11], [14], [33], [18] becomes even more challenging when the behaviors themselves are subtle. The adoption of explicit behaviors can be easily identified - for instance buying a DVD, joining a community, or using a hashtag in a tweet. Some behaviors are implicit, however, such as writing about of a topic, holding an opinion, or having a particular mood. In this paper, we focus on the diffusion of topics in social communities. Inferring topic diffusion introduces several additional challenges on top of the diffusion inference of explicit behaviors. First, topics are implicit and abstract concepts used in natural language. The adoption of topics cannot be directly identified, and instead has to be inferred from user-generated content. Second, the meaning of a topic evolves over time. A smart system should understand that '*MSN search*', '*Live search*', and '*Bing*' all refer to the same topic '*the Microsoft search engine*', with unique aspects at different time; and it should be able to track and adapt to this content change in the inference of topic diffusion. Third, information transmission is a complex social-psychological behavior [20], so the diffusion process of contents is inevitably influenced by the social relationships of the users. Moreover, the evolution and the diffusion of topics are compound processes: indeed, when a topic spreads from one user to another, new perspectives or new

focus is introduced to the topic; and an outbreak of a topic is usually accompanied by a shift of the meaning of the topic. Although there has been a line of work on the diffusion inference of explicit behaviors [6], [13], [9], [10], [16], [33], [25], [12], [1], [18], and a line of studies that incorporate network regularization into topic modeling [22], [4], [30], [5], none of this work addresses these challenges, making the existing methods incapable of accurately inferring the diffusion paths of topics.

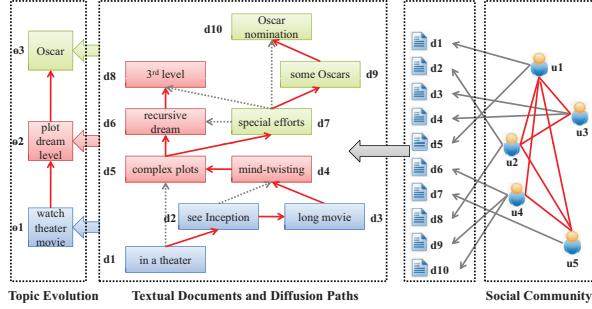


Figure 1. Example of Topic Diffusion and Evolution

In this paper, we address these challenges by studying the joint inference of topic diffusion and evolution in social communities. Content and linkage in user-generated text information, together with network structures, are used to facilitate the identification of topic adoption, the tracking of topic evolution, and the estimation of actual diffusion paths of any arbitrary topic (intuition illustrated in Figure 1).

When a topic is introduced into the community by a user, other users read the document(s) she wrote (e.g., tweets, blogs, scientific papers, etc) and adopt the topic by writing about it themselves. They may or may not cite the original document, or they may cite it together with other documents. Although topics are spread among documents rather than directly through social connections, we consider it much more likely for users to adopt ideas from their social connections (e.g., friends, people they follow, or people they have cited before) than from a stranger. Each document can not only adopt content from documents that influenced it, but may also include novel perspectives about the topic, and pass on the ‘innovation’ to other documents. The meaning of the topic thus evolves over time. The goal of the joint inference of topic diffusion and topic evolution is to identify the ‘real’ paths through which the topic propagates (red edges among documents in Figure 1), and also identify the specific temporal versions of the topic.

In this paper, we propose a novel statistical model for topic-based information diffusion and evolution (TIDE). Specifically, a mixture model is introduced to model the generation of text according to the diffusion and evolution of the topic, while the whole diffusion process is regularized with user-level social influences through a Gaussian Markov Random Field. The discovery of novel aspects and the

diffusion paths of the topic can be done by the joint inference of topic diffusion and evolution in TIDE.

## II. PROBLEM FORMULATION

In this section, we formally define the task of inferring the diffusion process and tracking topic evolutions in social communities. We begin with a few key concepts as follows.

**Definition II.1. Social Network.** A network is a graph  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ , where  $\mathcal{V}$  is a set of vertices and  $\mathcal{E}$  is a set of edges between vertices in  $\mathcal{V}$ . Particularly in a *social network*, a vertex corresponds to a user, and an edge  $e = (i, j)$  stands for a connection (or a tie) between two users  $i$  and  $j$ . The strength of the tie  $(i, j)$  is defined as a non-negative value  $g(i, j)$ . An edge can be either *directed* or *undirected*.

**Definition II.2. Document Collection.** A textual document  $d_i$  in a *document collection*  $\mathcal{D} = \{d_i\}_{i=1}^M$  is defined as a bag of words from a fixed vocabulary  $\mathcal{W} = \{w_k\}_{k=1}^L$ . That is,  $d_i = \{c(d_i, w_k)\}_{k=1}^L$ , where  $c(d, w)$  denotes the number of occurrences of word  $w$  in  $d$ .

**Definition II.3. Social Community.** A *social community* is defined as the union of a social network  $\mathcal{G}$  and a user-generated document collection  $\mathcal{D}$ , saying  $\{\mathcal{G}, \mathcal{D}\}$ . Each document  $d_i \in \mathcal{D}$  is associated with an *author*  $a_i$  in  $\mathcal{G}$  and a *time-stamps*  $t_i \in 1..T$ .

**Definition II.4. Topic.** A semantic *topic*  $\theta$  observed in a particular time period is defined as a multinomial distribution of words  $\{p(w|\theta)\}_{w \in \mathcal{W}}$  with the constraint  $\sum_{w \in \mathcal{W}} p(w|\theta) = 1$ .

**Definition II.5. Theme.** We define a general and coherent *theme* discussed in a social community as a stream of time-stamped topics  $\Theta = \{\theta_t\}_{t=0}^T$ . We call  $\theta_0$  the *primitive topic*, representing the original content of the theme prior to the discussions of the social community.  $\theta_{t>0}$  are time variant versions of  $\theta_0$ , which are gradually developed in the discussions of the social community, i.e.,  $\theta_t$  is the snapshot of  $\Theta$  at time  $t$ , which represents the novel aspect of the theme appearing at time  $t$ . Altogether  $\Theta$  represents the origin and evolution of the contents of the theme over time.

While the text content of individual documents can be explicitly observed, the general semantics of the time-variant topics and the adoption of the topic(s) in a document is implicit. Likewise, the source adopted in a document remains implicit. Naturally, there could be multiple sources: a document can be influenced by a few other documents, thus inheriting the topic from those documents. Some sources may be more influential than others. A document could also introduce original perspectives of the topic without being influenced by any existing document. The existence and strengths of document influences create the topic’s diffusion process, formally defined as a *diffusion graph*.

**Definition II.6. Diffusion Graph.** Given a theme  $\Theta$ , we define a *diffusion flow* from one document  $d_j$  to another  $d_i$  ( $t_j < t_i$ ) as the likelihood that  $d_i$  adopted the topic of  $\Theta$  due to the influence of  $d_j$ . The strength of such a diffusion flow is denoted as a positive value  $\pi_{i,j}$ . Note that  $d_i$  also introduces its novel perspective to  $\Theta$ . In this case, we assume there is a diffusion flow into  $d_i$  from the time-stamped topic  $\theta_{t_i}$ , with a strength  $\pi_{i,\theta}$ . Therefore, we define the *diffusion vector*  $\pi_i$  as the vector representing the strength of all the diffusion flows into  $d_i$ , i.e.,  $\pi(i) = \{\pi_{i,j}\}_{d_j \in \mathcal{D}} \cup \{\pi_{i,\theta}\}$ , with the constraint  $\sum_{d_j \in \mathcal{D}} \pi_{i,j} + \pi_{i,\theta} = 1$ . The union of diffusion flows into all documents in  $\mathcal{D}$  is the *diffusion graph*, i.e.,  $\Pi = \{\pi(i)\}_{d_i \in \mathcal{D}}$ . Clearly,  $\Pi$  is both weighted and directed.

Although the actual diffusion graph is unobserved, there are proxy networks that convey weaker signals in social communities. In many cases, a reference network (denoted as  $\mathcal{R}$ ) of the documents can be observed: a citation network of scientific publications, a hyperlink network of blog articles, or a tweet network of posters and followers. Intuitively, the diffusion network should be highly correlated with such a reference network  $\mathcal{R}$ . However, the actual diffusion network could still be substantially different from  $\mathcal{R}$ , because some influential references may be hidden, and some explicitly cited references are not actually influential.

Another signal of influence to consider is the social network structure. An author is likely to follow the work of his social connections, and thus is likely to adopt topics and ideas from the documents they generate [22], [30]. We refer to the set of documents pointing to  $d_i$  in the reference network as  $d_i$ 's *reference set*, denoted as  $r(i) \subset \mathcal{D}$ . When no signal of citation or social communication is available,  $r_i$  can be simply defined as all documents with a time stamp prior to  $t_i$ . When such a reference network is available, we assume  $\pi_{i,j} = 0$  if  $j \notin r(i)$ . Clearly, we also have  $\pi_{i,i} = 0$ .

Based on the above concept definitions, we can formalize the two major tasks of tracking **the diffusion and evolution of topics in social communities**. Given the input of a social community  $\mathcal{G}$ , a user-generated document collection  $\mathcal{D}$ , and the primitive topic  $\theta_0$  defining a theme, we aim to:

**Task 1: Infer the Diffusion Graph.** The goal is to discover the latent diffusion flow graph documents (and topics) (i.e.  $\Pi$ ). The result of this task can be used to discover (i) *the source(s) of topic in a document*: to what extent is the document influenced by other documents, and (ii) *the degree of originality in a document*: how much of a novel perspective does the document introduce to the topic.

**Task 2: Track Topic Evolution.** The goal is to infer the time-variant versions of topics (i.e.,  $\{\theta_t\}_{t=1}^T$ ) of a theme. By inferring  $\Theta$  given  $\theta_0$ , we expect to track new developments of the theme, understand its evolution over time, and better understand how it influences documents, etc.

### III. PROPOSED MODELS

In this section, we propose a novel and integrative probabilistic model of Text-based Information Diffusion and Evolution (TIDE) in social communities. Based on TIDE, we present the joint inference of the diffusion graph and the evolution of arbitrary topics.

#### A. Intuitions and the General Model

The general model of TIDE is designed based on a few key observations in social communities.

**Observation 1. Diffusion and Content.** When there is a significant diffusive flow between two documents, or one document significantly influences another, the contents of these two documents tend to be highly related. On the other hand, if two documents talk about different subjects, there is unlikely to be a salient influence flow or significant diffusion flow between them, even if one cites the other [27]. *W.l.o.g.*, we can assume that the content of a document depends on the documents that have influenced it.

**Observation 2. Diffusion and Social Connections.** Information transmission is a complex social-psychological behavior [20] (e.g., users tend to exhibit persistent interests [23]). The diffusion process among documents is likely to be regularized by the social connections of their authors. Indeed, an author is more likely to follow the works of her friends and would sooner adopt ideas from a friend than from a random author. The diffusion flows among documents are thus dependent to the social network of authors.

**Observation 3. Diffusion and Evolution.** Both the semantics of the topic and the regularization effect of the social network of users evolve over time. If an aspect in a document never appears in any of its potential references (in the form of either papers it cites or all existing papers available to its author(s)), it is likely that the aspect is an original idea, introduced by the document, which contributes to the evolution of the general theme. Meanwhile, the strength of influence through old social connections would decay after a reasonably long time.

Given a collection of authored and time-stamped documents  $\mathcal{D}$ , a social community  $\mathcal{G}$  of users who published these documents, and a primitive topic  $\theta_0$  representing the original semantics of a theme, we aim to inferring the latent stream of topics  $\Theta$  and the diffusion graph  $\Pi$ . Based on our observations above, the task of TIDE is then cast as the joint inference of the posterior of  $\Theta$  and  $\Pi$ . Formally, our objective is to infer:

$$P(\Pi, \Theta | \mathcal{G}, \mathcal{D}, \theta_0) \propto P(\Theta | \Pi, \mathcal{D}, \theta_0) \cdot P(\Pi | \mathcal{G}) \quad (1)$$

Based on our observations, we assume that the generation of the diffusion graph (only) depends on the social network structure, while the evolution of topics depends on the documents, the diffusion process, and of course the original version of the topic. We denote the first component of Equation 1 as the *topic model* and the second as the

*diffusion model*. Please note that although TIDE can be easily extended to model the mixture of multiple topics (similar to LDA [3]), we only present the primitive case to model one given topic. Our focus is to model the diffusion and evolution of any given topic rather than the discovery of multiple topics, which we leave to our future work.

In the topic model, a *mixture model* is designed to extract the topic snapshots (time-variant versions) of the theme (Section III-B). In the diffusion model, we introduce a *Gaussian markov random field* based on *graph projection* to model the dependency of diffusion flows on social connections (Section III-C). Finally, the inference of the combined model is discussed in Section III-D.

### B. The Topic Model

It is difficult to directly compute the posterior of topics  $\Theta$ . We make the following transformation:

$$P(\Theta|\Pi, \mathcal{D}, \theta_0) \propto P(\mathcal{D}|\Theta, \Pi, \theta_0) \cdot P(\Theta|\theta_0), \quad (2)$$

where the introduction of new aspects to the topic (i.e., the time-variant topic snapshots) does not depend on the diffusion flows.

We consider a typical generative process of  $\mathcal{D}$ : each document  $d_i$  is generated from a mixture model. When writing each word in  $d_i$ , one first chooses a component model from the mixture with a certain probability; once the component model  $\theta$  is selected, a word is sampled according to the word distribution of  $\theta$ .

We first introduce a background component model  $\theta_B$  estimated from the entire collection that explains the generation of common English words in the document  $d_i$ . The rest of the component models are designed based on the diffusion flows. Specifically, we introduce a component model for each document  $d_j$  that could have potentially influenced  $d_i$ . There is a non-trivial diffusion flow from  $d_j$  to  $d_i$ , and  $d_i$  could inherit the topic of  $d_j$  according to the strength of this diffusion. These component models can be estimated by simply using a maximum likelihood estimator on the corresponding  $d_j$ . Finally, we introduce a component model to explain the novel aspects introduced by the document  $d_i$ , i.e., the aspects that are not influenced by any existing documents. We assume that such aspects are generated directly from the latent topic at the time that  $d_i$  is written ( $\theta_{t_i}$ ). In other words, the original content is diffused from the topic directly to the document instead of from other documents. We assume that the probability of choosing each component is proportional to the strength of the diffusion vector, i.e.,  $\pi(i)$ .

Formally, the probability of generating a word  $w$  in  $d_i$  is  $p(w|d_i) = (1 - \lambda_B)(\sum_{j \in r(i)} \pi_{i,j} p(w|\theta_{d_j}) + \pi_{i,\theta} p(w|\theta_{t_i})) + \lambda_B p(w|\theta_B)$ , where  $\lambda_B$  is a predefined parameter that fixes the sampling probability of the background model. Note that for documents  $d_j \notin r(i)$ , we have  $\pi_{i,j} = 0$ . The likelihood

of the collection  $\mathcal{D}$  is given as:

$$P(\mathcal{D}|\Pi, \Theta, \theta_0) = \prod_{d_i \in \mathcal{D}} \prod_{w \in \mathcal{W}} p(w|d_i)^{c(w, d_i)}$$

We then consider the generation of the time-variant versions of the topic,  $\Theta$ . In TIDE, the primitive topic  $\theta_0$  is realized as a conjugate Dirichlet prior of the time-variant topic model  $\theta_t$ :  $Dir(\{1 + \mu_E p(w|\theta_0)\}_{w \in \mathcal{W}})$ . By doing so, we regularize these time-variant topic snapshots so that they can reflect the novel aspects of the theme, but do not shift away from it.  $\mu_E$  indicates how much we rely on the prior. Formally,

$$P(\Theta|\Pi) = \prod_{t \in 1..T} p(\theta_t|\theta_0) = \prod_{t \in 1..T} \prod_{w \in \mathcal{W}} p(w|\theta_t)^{\mu_E p(w|\theta_0)}$$

### C. The Diffusion Model

Compared to the modeling of topic evolution, the modeling of diffusion graph ( $P(\Pi|\mathcal{G})$ ) is less straightforward. Intuitively, the diffusion graph  $\Pi$  should be regularized by the social network  $\mathcal{G}$ , as social influence plays an important role in topic diffusion. However,  $\Pi$  is a network of *documents* while  $\mathcal{G}$  is a network of *users*. This makes it hard to model the regulation effect of  $\mathcal{G}$  on  $\Pi$ . We need a bridge between the two heterogeneous networks, for which we introduce the operation of *graph projection*.

**Definition III.1. Graph Projection.** Let  $\mathcal{G}_1$  and  $\mathcal{G}_2$  be two graphs. A projection  $f : \mathcal{G}_1|\mathcal{G}_2 \rightarrow \mathcal{G}'_1$  is called a *graph projection* if (i)  $\mathcal{V}(\mathcal{G}'_1) = \mathcal{V}(\mathcal{G}_2)$ ; (ii)  $\forall v \in \mathcal{V}(\mathcal{G}'_1), \exists u \in \mathcal{V}(\mathcal{G}_1)$  s.t.  $v \in f(u)$ ; and (iii)  $\forall e = (u, v) \in \mathcal{E}(\mathcal{G}'_1), \forall u' \in f(u)$  and  $v' \in f(v), e' = (u', v') \in \mathcal{E}(\mathcal{G}'_1)$ .

Through graph projection, the same vertex set is mapped into two networks, so that comparing the networks becomes more efficient and intuitive. Note that there are two asymmetric projection directions: 1) projecting  $\mathcal{G}$  into a document network and using it as an *a priori* of  $\Pi$ , or 2) projecting  $\Pi$  into a social network and considering the generation of such a social network based on  $\mathcal{G}$ . Since the document collection  $\mathcal{D}$  is usually much larger than the set of user  $\mathcal{V}(\mathcal{G})$ , projecting the document network into a social network has an unavoidable risk of losing information. Although this loss does not completely negate the value of the second direction of graph projection, in this work we consider the first direction: the projection of  $\mathcal{G}$  into a document network.

Denote  $\Pi'$  as the document network projected from  $\mathcal{G}$ , s.t.

$$P(\Pi|\mathcal{G}) = P(\Pi|\Pi') = P(\{\pi(i)\}_{d_i \in \mathcal{D}}|\Pi').$$

The remaining issue is how to fold  $\mathcal{G}$  into  $\Pi'$  and how to model the generation of  $\Pi$  based on  $\Pi'$ . Note that like  $\Pi$ , we can also denote  $\Pi' = \pi'(i)_{d_i \in \mathcal{D}}$ . We start with the generative model  $P(\Pi|\Pi')$ .

*Gaussian Graphical Models* (GGM) [34] are classical models used to explain the generation of networks, which could be an ideal solution to our problem. In a typical GGM, each nodes in the graph is modeled as a random

variable, for example a vector of  $k$  features. In our scenario, such a vector can be implemented as the diffusion vector  $\pi(i)$ . The joint distribution of all these variables (in our case,  $P(\pi(i))$ ) is assumed to be a multivariate Gaussian. Each edge in  $\Pi'$  represents the conditional dependency between two Gaussian variables, so the graph structure  $\Pi'$  corresponds to the inverse covariance matrix.

However, the computational complexity of such a graphical model usually scales cubically with the number of variables, and therefore becomes intolerable for even a moderately sized dataset. To make our model practical, we introduce an independency assumption: the diffusion vector of one document is independent of the others. By doing so, we can simplify the generative model of  $\Pi$  to be  $P(\Pi|\Pi') = \prod_{d_i \in \mathcal{D}} P(\pi(i)|\pi'(i))$ , where  $\pi'(i) = \{\pi'_{i,j}\}_{j \in r(i)} \cup \{\pi'_{i,\theta}\}$  is a conjugate prior vector, indicating the expected value of  $\pi(i)$ . Since  $\Pi'$  is projected from  $\mathcal{G}$ ,  $\pi'_{i,j}$  represents the social influence between  $a_j$  to  $a_i$ , which decays over time. By doing this, the document-level influence is regulated by the social tie at the user level.

Formally, we define  $\pi'_{i,j} = \frac{1}{Z(\pi'(i))} g(a_i, a_j) \cdot e^{-\frac{t_i - t_j}{\alpha}}$  by consolidating an exponential time model with  $\mathcal{G}$ . Intuitively, documents with higher authority are more likely to introduce more original content. We thus define  $\pi'_{i,\theta} = \frac{1}{Z(\pi'(i))} \text{Aut}(a_i)$ , where  $\text{Aut}(a_i)$  is an estimation of the authority of  $d_i$ .  $Z(\pi'(i))$  is a normalization factor such that  $\sum_{d_j \in \mathcal{D}} \pi'_{i,j} + \pi'_{i,\theta} = 1$ .

Given the design of  $\pi'(i)$ , the computation of  $P(\pi(i)|\pi'(i))$  is still non-trivial because of the dependency between the dimensions of  $\pi(i)$ . We introduce a *Gaussian Markov Random Field* [24] to model the conditional probability  $P(\pi(i)|\pi'(i))$  for each  $d_i$ .

**Definition III.2. Gaussian Markov Random Field (GMRF).** A random vector  $\xi = (x_1, x_2, \dots, x_n)^T$  is called a GMRF w.r.t. the graph  $\mathcal{G} = (\mathcal{V} = \{1, 2, \dots, n\}, \mathcal{E})$ , the mean  $\mu$  and the precision matrix  $\mathcal{Q}_\xi$ , iff the density of  $\xi$  s.t.,  $P(\xi) = (2\pi)^{-n/2} |\mathcal{Q}_\xi|^{1/2} e^{-\frac{1}{2}(\xi - \mu)^T \mathcal{Q}_\xi (\xi - \mu)}$  and  $\mathcal{Q}_\xi(i, j) \neq 0 \Leftrightarrow (i, j) \in \mathcal{E}$  for all  $i \neq j$ .

In our case, the random vector is the diffusion vector  $\pi(i)$ , with the mean as the prior vector  $\pi'(i)$ . The precision matrix  $\mathcal{Q}_{\pi(i)}$  corresponds to the similarities between the dimensions of  $\pi(i)$  (documents and topic snapshots), which can be expressed as the content similarities of corresponding  $\theta_{d_j}$ 's and  $\theta_t$ 's. Computationally,  $P(\pi(i)|\pi'(i))$  is defined as

$$P(\pi(i)|\pi'(i)) \propto e^{-\frac{1}{2} \sum_{i', j' \in \{r(i)\} \cup \{\theta\}} (\pi_{i,i'} - \mu_{i,i'}) \mathcal{Q}_{\pi(i)}(i', j') (\pi_{i,j'} - \mu_{i,j'})}$$

#### D. Parameter Estimation

Given our model defined above, we can fit the model to the data and estimate the parameters using a Maximum A Posterior estimator [29]. The Expectation Maximization (EM) algorithm [21] is applied, which iteratively computes a local maximum of the posterior. Computationally, the log

likelihood we want to maximize is:

$$\begin{aligned} E_{\Lambda^{(n-1)}} \{\log p(C|\Lambda)p(\Lambda)\} \propto & \sum_{d_i, w, d_j \in r(i)} c(d_i, w)(1 - z_{d_i, w}^{(n)}(\theta_B)) z_{d_i, w}^{(n)}(\theta_{d_j}) \log((1 - \lambda_B) \pi_{i,j} p(w|\theta_{d_j})) \\ & + \sum_{d_i, w} c(d_i, w)(1 - z_{d_i, w}^{(n)}(\theta_B)) z_{d_i, w}^{(n)}(\theta_{t_i}) \log((1 - \lambda_B) \pi_{i,E} p(w|\theta_{t_i})) \\ & + \sum_{d_i, w} c(d_i, w) z_{d_i, w}^{(n)}(\theta_B) \log(\lambda_B p(w|\theta_B)) + \mu_E \sum_{\theta_{t_i}, w} p(w|\theta_{t_i}) \log p(w|\theta_{t_i}) \\ & - \frac{\mu_G}{2} \sum_{d_i} \sum_{i', j' \in \mathcal{N}(i)} (\pi_{i,i'} - \mu_{i,i'}) \mathcal{Q}_{\pi(i)}(i', j') (\pi_{i,j'} - \mu_{i,j'}) \end{aligned} \quad (3)$$

Here  $\mu_G$  is a weight combining two components, and we use terms  $z_{d_i, w}(\cdot)$  instead of  $p(z_{d_i, w} = \cdot)$  to simplify notation

In the E-Step, we compute the expectation of the hidden variables:

$$\begin{aligned} z_{d_i, w}^{(n)}(\theta_{d_j}) &= \frac{\pi_{i,j}^{(n-1)} p(w|\theta_{d_j})}{\sum_{j' \in r(i)} \pi_{i,j'}^{(n-1)} p(w|\theta_{d_{j'}}) + \pi_{i,\theta}^{(n-1)} p(w|\theta_{t_i})} \\ z_{d_i, w}^{(n)}(\theta_{t_i}) &= \frac{\pi_{i,\theta}^{(n-1)} p(w|\theta_{t_i})}{\sum_{j' \in r(i)} \pi_{i,j'}^{(n-1)} p(w|\theta_{d_{j'}}) + \pi_{i,\theta}^{(n-1)} p(w|\theta_{t_i})} \\ z_{d_i, w}^{(n)}(\theta_B) &= \frac{\lambda_B p(w|\theta_B)}{(1 - \lambda_B) (\sum_{j' \in r(i)} \pi_{i,j'}^{(n-1)} p(w|\theta_{d_{j'}}) + \pi_{i,\theta}^{(n-1)} p(w|\theta_{t_i})) + \lambda_B p(w|\theta_B)} \end{aligned}$$

In the M-step, given the expectation of the hidden variables, we get the best parameters  $p(w|\theta_t)$  as:

$$p(w|\theta_t) = \frac{\sum_{d_i, t_i=t} c(d_i, w)(1 - z_{d_i, w}^{(n)}(\theta_B)) z_{d_i, w}^{(n)}(\theta_t) + \mu_E p(w|\theta_0)}{\sum_{w, d_i, t_i=t} c(d_i, w)(1 - z_{d_i, w}^{(n)}(\theta_B)) z_{d_i, w}^{(n)}(\theta_t) + \mu_E p(w|\theta_0)}$$

By integrating Lagrange multipliers [21]  $f_i$  for each  $d_i \in \mathcal{D}$ , the inference of  $\pi(i)$  boils down to solving a group of cubic equations:  $\pi_{i,*}^2 + \beta_{i,*} \pi_{i,*} + \gamma_{i,*} = 0$  ( $*$   $\in r(i) \cup \{\theta\}$ ), where  $\beta_{i,*} = \frac{\sum_{* \neq *'} (\mathcal{Q}_{\pi(i)}(*, *') + \mathcal{Q}_{\pi(i)}(*', *)) (\pi_{i,*'}^{(n-1)} - \mu_{i,*'})}{2\mathcal{Q}_{\pi(i)}(*, *)} - \mu_{i,*} + \frac{f_i}{\mu_G \mathcal{Q}(i)_{*,*}}$  and  $\gamma_{i,*} = -\frac{\sum_w c(d_i, w)(1 - z_{d_i, w}^{(n)}(\theta_B)) z_{d_i, w}^{(n)}(\theta_{d_j})}{\mu_G \mathcal{Q}(i)_{*,*}}$ .

It is easy to prove that there exist valid solutions for the group of equations that satisfy the constraint  $\sum_{* \in r(i) \cup \{\theta\}} \pi_{i,*} = 1$  for each  $d_i$  in  $\mathcal{D}$ .

## IV. EXPERIMENTS

In this section, we evaluate the effectiveness of our TIDE model on synthetic datasets, and data collected from two real-world social communities, *DBLP* [32] and *Twitter* [19].

### A. Experimental Setup

1) *Data Collection: The DBLP Dataset* ([32]). The Digital Bibliography and Library Project (DBLP) is a web accessible database of the bibliographic information of computer science publications. In this experiment, we

use a collection of DBLP articles augmented with citation information, released by the ArnetMiner group, which contains 1,632,442 publications by 1,741,170 researchers with 2,327,450 citations. After filtering out papers without text or citation information, 243,425 papers and 246,839 authors remain. This dataset represents a typical academic community, with a social network of authors (with coauthor and citation relations) and a collection of scientific papers.

**The Twitter Dataset** ([19]). Twitter is a well known social networking and microblogging community. In this experiment, the dataset was crawled by the DAIS group at the University of Illinois, which contains 5,000 socially connected users and their most recent 200 tweets posted before Nov. 23, 2010. In total, there are 103,968 one-way following relations, and 51,032 pairs of Twitter friends (defined as a mutual following relation). This dataset represents a typical social community with a directed social network (defined by following relations) and a collection of tweets.

**Synthetic Dataset.** The lack of ground truth on a real world dataset makes it hard to evaluate the model performance quantitatively. For quantitative evaluation, we construct a synthetic dataset which simulates the diffusion of 1,000 themes. For each theme, we extract a subgraph of 1,000 authors from the *DBLP* dataset using breadth first search from a random seed author. This subgraph is used to simulate the social network in which the theme diffuses. We then randomly attach 1,859 empty and time-stamped documents to the authors in this network<sup>1</sup>. We then simulate a diffusion graph of the 1,859 documents that is regularized by the simulated social network structure. Specifically, we first randomly generate a network of the 1,859 documents using an Erdos/Renyi model, with an average degree of 5 (consistent with actual statistics in the *DBLP* dataset). The direction of each edge is determined by the time stamps of the documents (always pointing from an ‘older’ document to a ‘newer’ document). We then assign a weight to each edge based on the social connections of the authors of the two document, plus a random effect. This directed and re-weighted random network simulates a real diffusion network among documents. For each theme, we also simulate a sequence of 10 evolving topic snapshots based on the dynamic topic models [2]. Finally, the text content of each document is generated by a simple mixture model with all documents that have ‘influenced’ this document as well as the corresponding topic snapshot.

2) **Baselines: The NetInf Model** [9]. NetInf is a typical model that infers the diffusion network of explicit user behaviors. Given the time stamps at which individuals adopt a behavior, *NetInf* identifies the optimal general network of users that best explains the observed adoptions. Compared to *TIDE*, *NetInf* infers the general social network structure by observing the propagations of a group of events, while

<sup>1</sup>According to the statistics on our *DBLP* dataset, each researcher has 1.859 first-authored publications in average.

*TIDE* infers the theme-specific diffusion graph with the help of a general social network. Note *NetInf* does not consider text information, and thus cannot track topic evolution.

If we treat each term with a positive probability in the primitive topic as an explicit event/behavior, then a document adopts that behavior explicitly if the term appears in the document. We are then able to infer the optimal document network using *NetInf*. This optimal network is easily converted into a diffusion graph by endowing each edge with equal flow volume.

**The IndCas Model** [25]. The second baseline is a deviation of the independent cascade model stated in [25], where the probability for an active document to infect another is proportional to the strength of the social connection between their authors with an exponential decay effect [7] (see Section III-C). We convert these probabilities into a diffusion graph where the diffusion flow from  $d_j$  to  $d_i$  is proportional to the probability that  $d_j$  infects  $d_i$ .

**The TIDE— Model.** To evaluate the effectiveness of social connections in our models, we implement a special version of *TIDE* by removing the regularization term with the network structure, i.e., by setting  $\mu_G = 0$ .

We believe *NetInf* and *IndCas* are good representatives of diffusion inference models of explicit behaviors, which do not consider textual information or topic evolutions. *TIDE*— on the other hand ignores the effects of social connections.

## B. Experiments on Synthetic Data

The goal of the experiments on synthetic data is to quantitatively evaluate how well each method can (i) infer diffusion graphs, (ii) estimate the novelty of the contribution, and (iii) discover snapshots in topic evolution (if possible). Given the simulated social community (the social network, the document collection, and the primitive topic), our goal is to recover the diffusion graph and the topic snapshots. The parameters in the *TIDE* model are set empirically as  $\mu_E = 10$ ,  $\alpha = 30$ , and  $\mu_G = 10$ .

1) **Analysis on Information Diffusion:** Let us first introduce the evaluation metrics.

**Definition IV.1. Graph KL-Divergence.** The symmetrized Kullback-Leibler divergence [15] is a classic measure of the difference between two probability distributions. We extend the SKLD and define an evaluation metric to measure the discrepancy between two diffusion graphs  $\Pi_{\mathcal{P}}$  and  $\Pi_{\mathcal{Q}}$  on the same document collection  $\mathcal{D}$  as  $GD_{KL}(\Pi_{\mathcal{P}}, \Pi_{\mathcal{Q}}) = \frac{\sum_{d_i \in \mathcal{D}} (D_{KL}(\pi_{\mathcal{P}}(i) || \pi_{\mathcal{Q}}(i)) + D_{KL}(\pi_{\mathcal{Q}}(i) || \pi_{\mathcal{P}}(i)))}{2|\mathcal{D}|}$ .

**Definition IV.2. Graph Cosine Similarity.** We define a metric of similarity between two diffusion graphs  $\Pi_{\mathcal{P}}$  and  $\Pi_{\mathcal{Q}}$ , as the average cosine similarity [31] between their diffusion vectors as  $Cos(\Pi_{\mathcal{P}}, \Pi_{\mathcal{Q}}) = \frac{1}{|\mathcal{D}|} \sum_{d_i \in \mathcal{D}} \frac{\pi_{\mathcal{P}}(i) \cdot \pi_{\mathcal{Q}}(i)}{||\pi_{\mathcal{P}}(i)|| \cdot ||\pi_{\mathcal{Q}}(i)||}$ .

A better model should infer a diffusion graph that is closer

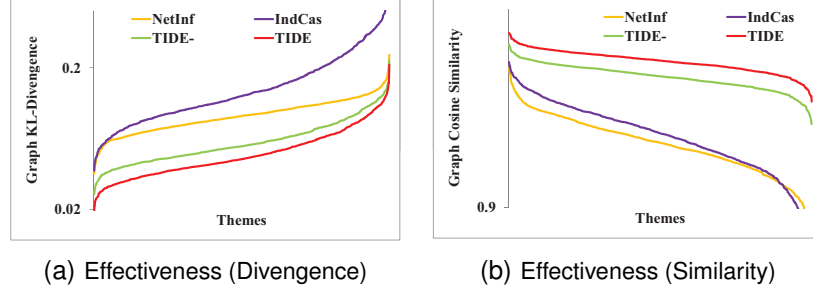


Figure 2. Diffusion Evaluation on the Synthetic Dataset

to the ‘ground truth’, that is, having a lower KL-divergence and a higher Cosine similarity.

Object	TDG		$\mathcal{H}$		$\mathcal{G}$	
Metric	GKLD	GCS	GKLD	GCS	GKLD	GCS
NetInf	0.094	0.936	1.291	0.869	0.949	0.802
IndCas	0.160	0.931	1.297	0.855	<b>0.721</b>	<b>0.898</b>
TIDE–	0.063*	0.969*	<b>0.991</b>	<b>0.949</b>	0.876	0.841
TIDE	<b>0.052*</b>	<b>0.972*</b>	1.011	0.938	0.846	0.855

Table I

DIFFUSION EVALUATION ON THE SYNTHETIC DATASET (TDG = TRUE DIFFUSION GRAPH, GKLD = GRAPH KULLBACK-LEIBLER DIVERGENCE, GCS = GRAPH COSINE SIMILARITY)

We measure the statistical significance of the improvement using the dependent t-test. \* means that the improvement (over the row above) hypothesis is accepted at significance level 0.001.

In practice, we calculate the two metrics for the result of each method and each theme, and connect the KL-divergence scores in decreasing order (Figure 2(a)) and Cosine-similarity scores in increasing order (Figure 2(b)). The aggregated performance of the 1,000 themes is reported in the 1<sup>st</sup> and 2<sup>nd</sup> columns of Table I. **We conclude that TIDE achieves the best performance, followed by TIDE–, NetInf, and IndCas.**

2) *Proof of Combined Power:* With this experiment, we can also prove that both social networks and text information play an important role the inference of topic diffusion. First, we create a document network (denoted as  $\mathcal{H}$ ), where the edge weight is proportional to the content similarities between documents. We compare each inferred diffusion graph with  $\mathcal{H}$ , and report the aggregated value of the two metrics in the 3<sup>rd</sup> and 4<sup>th</sup> columns of Table I. Second, we project each diffusion graph  $\Pi$  into a user network (denoted as  $f(\Pi)$ ), compare  $f(\Pi)$  with the general social network  $\mathcal{G}$ , and report the aggregated value of the two metrics in the last two columns of Table I. We can observe some phenomena that accord with our hypothesis in designing our model: TIDE– infers diffusion graphs only considering textual information without considering the social network structure, while IndCas infers the diffusion network purely

based on the social influences. Indeed, the diffusion networks inferred by TIDE– are significantly biased towards the document similarity networks  $\mathcal{H}$ , and the diffusion networks inferred by IndCas are biased towards the social networks  $\mathcal{G}$ . Neither of them infers diffusion networks that are closer to the ground truth than TIDE, which employs both text information and the social network.

Metric	KLD	CS
FM	0.4281	0.7033
TIDE–	0.3301*	0.8622*
TIDE	<b>0.2893*</b>	<b>0.8774*</b>

Table IV

EVOLUTION EVALUATION ON THE SYNTHETIC DATASET (KLD = KULLBACK-LEIBLER DIVERGENCE, CS = COSINE SIMILARITY, FM = FEEDBACK MODEL [36])

\* means the improvement (over the above row) hypothesis is accepted at the significance level 0.001 based on dependent t-test.

3) *Analysis on Content Evolution:* In this experiment, we study how successfully TIDE and the baseline models track topic evolution. Since NetInf and IndCas are not able to handle topics, we compare our models TIDE and TIDE– with a simple mixture model stated in [36].

We replicate the experiments described in Section IV-B1. We use two similar metrics (i.e., the symmetrized KL-divergence and the Cosine similarity) to measure the closeness of the discovered word distributions of the topic snapshots to the “ground truth” (topic snapshots we construct in the synthetic dataset). The results are reported in Table IV.

As shown in Table IV, TIDE outperforms the other two methods by a significant margin, which proves our statement in Section I: the evolution and the diffusion of topics are compound processes, and the success of one aspect will help the inference of the other.

### C. Experiments on Real Social Networks

We present the experiments on real world social communities in this section. Note that ‘ground truth’ diffusion networks and topic snapshots are usually not available.



ID	Publication	ID	Publication	ID	Tweet (incomplete)
A	J. Han, SIGMOD'00.	B	A. Khan, KDD'10.	A	Inception had better special effects than Videodrome.
C	X. Yan, SIGMOD'04.	D	M. Zaki, KDD'03.	B	Inception's effects might take some Oscars.
E	X. Yan, KDD'03.	F	Y. Chi, TKDE'05.	C	I predict Inception's 12 Oscar nominations.
G	M. Zaki, KDD'02.	H	A. Bifet, KDD'08.	D	It has to be like a 3rd level Inception dream.
I	X. Yan, KDD'05.	J	U. R��kert, SAC'04.	E	I wonder what level of recursive dreams.
K	C. Chen, CIKM'08.	L	J. Wang, KDD'03.	F	Inception. What a brilliant, mind-twisting movie.
M	J. Wang, TKDE'05.	N	F. Pan, KDD'03	G	Watching inception. Long movie.
O	A. Lee, Infomation System'10.			H	You'd be odd on twitter if you haven't seen Inception.
P	U. Yun, Knowledge-Based System'08			I	First time I have seen a movie in a theater in the last 6 months.
Q	J. Balc��zar, Machine Learning'10.			J	If you like intelligent movies and complex plots, see Inception.

Table II  
PUBLICATIONS SHOWN IN FIGURE 4(A)

Table III  
TWEETS SHOWN IN FIGURE 4(B)

1) *Verifying Motivating Observations*: We start with the verification of the authenticity of the three motivating observations stated in Section III-A. We expect social influence to play a role, so that an author is more likely to adopt topics from her social connections' documents. If this is the case, an author will consistently pay attention to papers published by authors that she knows, or has cited before. One intuitive way to verify this is through the behavior of 're-citation', i.e., once the author cited one paper, it is likely that she will cite a paper by the same author again. We group authors by the number of publications, and plot the average ratio of re-citations in Figure 3(a). It shows that there are substantial re-citation behaviors, and as an author publishes more papers, the ratio of re-citation also grows. This verifies the existence of social influence in document-level information diffusion.

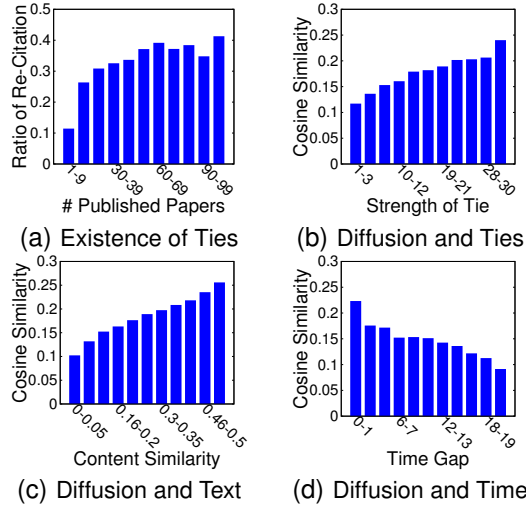


Figure 3. Verifying Observations by DBLP-Citation Dataset

Rather than inferring the diffusion, a rough proxy of the influence of a cited paper  $d_r$  on a citing paper  $d_c$  can be measured by the authors' phrasing after the citation.

Generally, if the authors of  $d_c$  publish many papers related to  $d_r$  after they publish  $d_c$ , it is fair to believe  $d_r$  is quite influential on  $d_c$ . We partition the citations (each of which is defined by a cited paper  $d_r$  and a citing paper  $d_c$ ) into different groups according to the strength of the social connection between their authors. For each citation, we then compute the average document similarity between  $d_r$  and all papers published by  $d_c$ 's authors after they had published  $d_c$ . The aggregated similarity is plotted in Figure 3(b). We repeat the same experiment, but partition citations by degrees of the content similarity of  $d_r$  and  $d_c$  (Figure 3(c)), as well as the time gap (3(d)) between  $d_r$  and  $d_c$ . Figure 3(b)-3(d) support our hypotheses that the (proxy) influence between two documents increases with the strength of social ties (Observation 2) and the content similarity, but decays over time (Observation 3).

2) *Case Study*: We select two themes for case study: one is about the research topic 'frequent pattern mining' in the DBLP-Citation dataset (see the 1<sup>st</sup> column of Table V), and the other is about the movie 'Inception' in the Twitter dataset (see the 1<sup>st</sup> column of Table VII).

**Analysis on Information Diffusion.** For theme 1, we apply the *TIDE* model to 344 papers published during the past ten years (2000 to 2010), which contain at least three primitive keywords in the title or abstract. A subgraph of the diffusion graph estimated by TIDE is shown in Figure 4(a), on a subset of 17 selected papers (listed in Table II). The volume of each diffusion flow is marked on the edge. To quantitatively access the result, we compare the graph with three alternative "diffusion graphs." In the first graph, the weight of an edge  $d_r \rightarrow d_c$  is set to be proportional to the total length of citation sentences where  $d_c$  mentions  $d_r$ . We then employ two experts to manually score the impact of each reference paper in a scale from one to five. The **Mean Absolute Error** [26], as the statistical metric of accuracy, based on each criteria, and the **Cohen's Kappa Coefficient** [28], as the



Primitive Topic		Year 2003		Year 2005		Year 2009		Year 2003		Year 2005		Year 2009	
frequent	0.20	itemset	0.05	itemset	0.04	itemset	0.03	efficient	0.02	close	0.01	sequential	0.02
pattern	0.40	GSM	0.03	tree	0.02	tree	0.02	close	0.01	itemset	0.01	itemset	0.01
mining	0.20	association	0.02	parallel	0.01	sequence	0.01	association	0.01	match	0.01	tree	0.01
graph	0.05	apriori	0.02	graph	0.01	graph	0.01	support	0.01	tree	0.01	graph	0.01
tree	0.05	tree	0.01	sequence	0.01	slide	0.01	query	0.01	graph	0.01	database	0.01
sequence	0.05	graph	0.01	traversal	0.01	gram	0.01	temporal	0.01	sequential	0.01	efficient	0.01
itemset	0.05	subgroup	0.01	optimize	0.01	window	0.01	graph	0.01	efficient	0.01	rule	0.01
		sequential	0.01	suffix	0.01	apriori	0.01	rule	0.01	application	0.01	match	0.01

Table V  
TOPIC SNAPSHOTS BY *TIDE* ON THEME 1 (DBLP)

Table VI  
TOPIC SNAPSHOTS BY [36] ON THEME 1 (DBLP)

Primitive Topic		Jul 16-19		Jul 20-23		Jul 24-27		Jul 16-19		Jul 20-23		Jul 24-27	
inception	1.00	watch	0.05	dream	0.06	oscar	0.04	movie	0.06	type	0.05	oscar	0.03
		night	0.05	mind	0.05	effect	0.04	night	0.06	eye	0.05	act	0.03
		movie	0.05	level	0.03	dream	0.02	special	0.03	watch	0.05	dream	0.03
		special	0.03	walk	0.01	clever	0.01	watch	0.03	night	0.05	strong	0.02
		enjoy	0.01	recursive	0.01	briliant	0.01	bad	0.03	dream	0.04	night	0.02

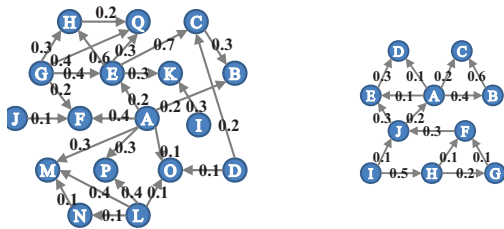
Table VII  
TOPIC SNAPSHOTS BY *TIDE* ON THEME 2 (TWITTER)

Table VIII  
TOPIC SNAPSHOTS BY [36] ON THEME 2 (TWITTER)

Table IX  
CASE STUDY ON REAL NETWORKS: TOPIC EVOLUTION

measure of inter-criteria agreement, are reported in Table X.

Theme 2 has been used as the running example in Section I (see Figure 1), and allows us to reveal more details. We apply the *TIDE* model on 361 tweets containing the keyword ‘*Inception*’, and draw the diffusion graph on 10 selected tweets (listed in Table III) in Figure 4(b). We repeat the same evaluation procedure as was done for theme 1 (see Table XI), but the edge weight of the first criteria graph is based on whether one tweet was replying to the other.



(a) Theme 1

(b) Theme 2

Figure 4. Case Study on Real Networks: Diffusion Graphs

In both cases, the opinions of the first expert gains the most agreement from others, and our result has the highest accuracy against the ground truth given by the first expert.

**Analysis on Content Evolution.** We apply both *TIDE* and the feedback model [36] to extract topic snapshots for two themes. Top words (with probabilities) of several selected topics are listed in Table V-VIII. To show more results, the word ‘*frequent*’, ‘*pattern*’ and ‘*mining*’ are eliminated from Table V and VI; and the word ‘*Inception*’ is eliminated from

MAE	SL	Exp1	Exp2	CKC	Exp1	Exp2
TIDE	0.122	<b>0.108</b>	0.120	SL	0.502	0.210
				Exp1	–	<b>0.633</b>

Table X  
THEME 1 (DBLP)

MAE	RR	Exp1	Exp2	CKC	Exp1	Exp2
TIDE	0.363	<b>0.130</b>	0.135	RR	0.358	0.373
				Exp1	–	<b>0.750</b>

Table XI  
THEME 2 (TWITTER)

Table XII  
CASE STUDY ON REAL NETWORKS: ACCURACY EVALUATION (MAE = MEAN ABSOLUTE ERROR, CKC = COHEN’S KAPPA COEFFICIENCY, SL = SENTENCE LENGTH, RR = REPLYING RELATION)

Table VII and VIII.

As described in Section III-B, the topic component of *TIDE* utilize (i) a background model to absorb common words, and (ii) reference models to absorb old words, so that topic snapshots would attract more discriminative and meaningful words that describe the novel aspect of a theme. For example, the topic at ‘*Year 2009*’ in Table VI reveals a new trend of mining patterns up to a certain length (i.e. ‘*gram*’) in a ‘*sliding*’ ‘*window*’, and the topic at ‘*Jul 20-23*’ in Table VIII talks about the movie plots such as the ‘*level*’ of a ‘*recursive*’ ‘*dream*’. However, since [36] only considers the idea of a background model, these interesting

new words are easily overlooked, because antiquated words, such as ‘efficient’ in Table V and ‘watch’ in Table VII, repeatedly appear in lots of topic snapshots.

## V. CONCLUSION

In this paper, we propose *TIDE*, a novel probabilistic model for the joint inference of diffusion and evolution of topics in social communities. *TIDE* integrates the generation of text, the evolution of topics, and the social network structure in a unified model. Given the primitive form of any arbitrary topic, *TIDE* effectively tracks the topic snapshots that evolve over time and reveals the latent diffusion paths of the topic. Comprehensive experimental studies on both synthetic data and two real-world datasets show that *TIDE* outperforms existing approaches.

One important finding is that the discovery of topic diffusion and topic evolution benefits significantly from the joint inference process. Text information, social influence, and the general social network structure play very important roles in the inference process. We plan to develop a future extension of *TIDE* that models the evolution of the social network structure in addition to the evolution of topics.

## ACKNOWLEDGEMENT

This work at the University of Illinois was supported in part by the Army Research Laboratory accomplished under Cooperative Agreement Number W911NF-09-2-0053, NSF grant IIS-09-05215, and NASA grant NNX08AC35A. The views and conclusions contained in this document are those of the authors and should not be interpreted as representing the official policies of the U.S. Government. The U.S. Government is authorized to reproduce and distribute reprints for Government purposes notwithstanding any copyright notation stated here on. The work at the University of Michigan is support by the national science foundation under grant number IIS-0968489 and IIS-1054199.

## REFERENCES

- [1] L. Backstrom, D. P. Huttenlocher, J. M. Kleinberg, and X. Lan. Group formation in large social networks: membership, growth, and evolution. In *KDD*, 2006.
- [2] D. M. Blei and J. D. Lafferty. Dynamic topic models. In *ICML*, 2006.
- [3] D. M. Blei, A. Y. Ng, and M. I. Jordan. Latent dirichlet allocation. In *NIPS*, 2001.
- [4] D. Cai, Q. Mei, J. Han, and C. Zhai. Modeling hidden topics on document manifold. In *CIKM*, 2008.
- [5] D. Cai, X. Wang, and X. He. Probabilistic dyadic data analysis with local and global consistency. In *ICML*, 2009.
- [6] W. Chen, C. Wang, and Y. Wang. Scalable influence maximization for prevalent viral marketing in large-scale social networks. In *KDD*, 2010.
- [7] R. Crane and D. Sornette. Robust dynamic classes revealed by measuring the response function of a social system. *PNAS*, 2008.
- [8] D. Easley and J. Kleinberg. Networks, crowds, and markets: Reasoning about a highly connected world. *Cambridge University Press*, 2010.
- [9] M. Gomez-Rodriguez, J. Leskovec, and A. Krause. Inferring networks of diffusion and influence. In *KDD*, 2010.
- [10] A. Goyal, F. Bonchi, and L. V. S. Lakshmanan. Learning influence probabilities in social networks. In *WSDM*, 2010.
- [11] D. Gruhl, R. V. Guha, D. Liben-Nowell, and A. Tomkins. Information diffusion through blogspace. In *WWW*, 2004.
- [12] M. Kimura, K. Saito, and H. Motoda. Minimizing the spread of contamination by blocking links in a network. In *AAAI*, 2008.
- [13] M. Kimura, K. Saito, R. Nakano, and H. Motoda. Extracting influential nodes on a social network for information diffusion. In *Data Min. Knowl. Discov.*, 2010.
- [14] M. Kimura, K. Saito, K. Ohara, and H. Motoda. Learning to predict opinion share in social networks. In *AAAI*, 2010.
- [15] S. Kullback and R. Leibler. On information and sufficiency. *Annals of Mathematical Statistics*, 1951.
- [16] C. Lee, H. Kwak, H. Park, and S. B. Moon. Finding influentials based on the temporal order of information adoption in twitter. In *WWW*, 2010.
- [17] J. Leskovec, L. A. Adamic, and B. A. Huberman. The dynamics of viral marketing. *TWEB*, 2007.
- [18] J. Leskovec, M. McGlohon, C. Faloutsos, N. S. Glance, and M. Hurst. Patterns of cascading behavior in large blog graphs. In *SDM*, 2007.
- [19] C. X. Lin, B. Zhao, Q. Mei, and J. Han. Pet: a statistical model for popular events tracking in social communities. In *KDD*, 2010.
- [20] M. H. MacRobert and B. R. MacRoberts. Problems of citation analysis. *Scientometrics*, 1996.
- [21] G. McLachlan and T. Krishnan. The EM algorithm and extensions. *Wiley Series in Probability and Statistics*, 2001.
- [22] Q. Mei, D. Cai, D. Zhang, and C. Zhai. Topic modeling with network regularization. In *WWW*, 2008.
- [23] F. Qiu and J. Cho. Automatic identification of user interest for personalized search. In *WWW*, 2006.
- [24] H. Rue and L. Held. Gaussian markov random fields: Theory and applications - theory and application. *Chapman and Hall/CRC*, 2006.
- [25] K. Saito, M. Kimura, K. Ohara, and H. Motoda. Selecting information diffusion models over social networks for behavioral analysis. In *ECML/PKDD*, 2010.
- [26] B. M. Sarwar, G. Karypis, J. A. Konstan, and J. Riedl. Item-based collaborative filtering recommendation algorithms. In *WWW*, 2001.
- [27] A. Si, H. V. Leong, and R. W. H. Lau. Check: a document plagiarism detection system. In *SAC*, 1997.
- [28] N. Smeeton. Early history of the kappa statistic. *Biometrics*, 1985.
- [29] H. W. Sorenson. Parameter estimation: Principles and problems. *M. Dekker*, 1980.
- [30] Y. Sun, J. Han, J. Gao, and Y. Yu. itopicmodel: Information network-integrated topic modeling. In *ICDM*, 2009.
- [31] P. Tan, M. Steinbach, and V. Kumar. Introduction to data mining. *Addison Wesley*, 2005.
- [32] J. Tang, J. Zhang, L. Yao, J. Li, L. Zhang, and Z. Su. Arnetminer: extraction and mining of academic social networks. In *KDD*, 2008.
- [33] X. Wan and J. Yang. Learning information diffusion process on the web. In *WWW*, 2007.
- [34] J. Whittaker. Graphical models in applied multivariate statistics. *Wiley Publishing*, 2009.
- [35] J. Yang and J. Leskovec. Modeling information diffusion in implicit networks. In *ICDM*, 2010.
- [36] C. Zhai and J. D. Lafferty. Model-based feedback in the language modeling approach to information retrieval. In *CIKM*, 2001.