



# Nonnegative matrix factorization algorithms for link prediction in temporal networks using graph communicability



Xiaoke Ma<sup>a,1,\*</sup>, Penggang Sun<sup>a,1</sup>, Guimin Qin<sup>b,1</sup>

<sup>a</sup>School of Computer Science and Technology, Xidian University, No.2 South Taibai Road, Xi'an, Shaanxi, China

<sup>b</sup>School of Software, Xidian University, No.2 South Taibai Road, Xi'an, Shaanxi, China

## ARTICLE INFO

### Article history:

Received 27 September 2016

Revised 26 May 2017

Accepted 16 June 2017

Available online 17 June 2017

### Keywords:

Dynamic networks

Nonnegative matrix factorization

Eigenvalues and eigenvector

Temporal link prediction

## ABSTRACT

Networks derived from many disciplines, such as social relations, web contents, and cancer progression, are temporal and incomplete. Link prediction in temporal networks is of theoretical interest and practical significance because spurious links are critical for investigating evolving mechanisms. In this study, we address the temporal link prediction problem in networks, i.e. predicting links at time  $T + 1$  based on a given temporal network from time 1 to  $T$ . To address the relationships among matrix decomposition-based algorithms, we prove the equivalence between the eigendecomposition and nonnegative matrix factorization (NMF) algorithms, which serves as the theoretical foundation for designing NMF-based algorithms for temporal link prediction. A novel NMF-based algorithm is proposed based on such equivalence. The algorithm factorizes each network to obtain features using graph communicability, and then collapses the feature matrices to predict temporal links. Compared with state-of-the-art methods, the proposed algorithm exhibits significantly improved accuracy by avoiding the collapse of temporal networks. Experimental results of a number of artificial and real temporal networks illustrate that the proposed method is not only more accurate but also more robust than state-of-the-art approaches.

© 2017 Elsevier Ltd. All rights reserved.

## 1. Introduction

The network (sometimes called graph) effectively characterizes and analyzes complex systems, in which each vertex represents an individual, such as a biological entity (e.g., a gene or a protein), a web user, or a terminal in Internet. Each link denotes an interaction between a pair of vertices. Various real-world networks have been derived from, such as social networks [1,2], technological networks [3] and biological networks [4]. Network analysis has emerged as a key technique in modern science with the immediate purpose of discovering graph patterns by elucidating the structure-function relationship of overall systems. For example, communities in protein interaction networks correspond to the protein complexes that are critical for biological processes [5].

However, many networks are incomplete because of the limitations of our knowledge regarding complex systems, which significantly hinder the practical application of network analysis. For example, nearly 80% of interactions within yeast [6] and 99% within human [7] remain unknown. Accordingly, link prediction plays a critical role in network analysis [8,9], which does not only help us

recover data, but also improve our understanding of the mechanisms of networks.

Therefore, approaches for predicting links in networks are urgently required, and considerable efforts have been exerted to address this issue [9–12]. Available methods for link prediction can be categorized into two classes: experimental and computational methods. Experimental methods use a physical strategy to validate the existence of links. They fail to provide satisfactory answers primarily due to the limitations of finance and technology. These methods are also costly and time-consuming, particularly for validating interactions among proteins through biological experiments [13]. Thus, computational methods for predicting links based on known interactions becomes alternatives for experimental approaches [14–19].

However, the vast majority of available algorithms focus on static networks, where the ultimate goal is to predict links to describe a complete picture of the whole network structure [3]. Many networks derived from the real world dynamically change over time (called temporal or dynamic networks) [20]. For example, in scientific collaboration networks, interactions evolve since scientists directly change their collaborators as they shift their research directions [21]. In disease networks, cancer metastasis is mainly due to cancer cell immigration [22]. Thus, the analysis of temporal networks has received considerable attention because the evo-

\* Corresponding author.

E-mail addresses: [xkma@xidian.edu.cn](mailto:xkma@xidian.edu.cn), [maxiaoke8218@163.com](mailto:maxiaoke8218@163.com) (X. Ma).

<sup>1</sup> Equal contribution.

lution patterns provide novel insights into the underlying mechanisms of complex networks [20,23–26].

Accordingly, the prediction of links in temporal networks is a promising and interesting subject because it is the foundation for network analysis. Unlike *missing link prediction* in static networks, the *temporal link prediction problem* obtains the edges in a network at time  $T + 1$  based on given temporal network from time 1 to  $T$ . This problem is applied to a variety of contexts, such as collaborative filtering [27] and social network connections [28]. However, designing algorithms for the temporal link prediction problem is highly non-trivial [3] because of two reasons. First, features in temporal networks are significantly more complicated than those in static ones, and thus, they are difficult to characterize and extract. Second, the complexity of temporal networks poses a considerable challenge to designing effective and efficient algorithms.

Although the process is difficult, many algorithms for predicting temporal links have been proposed [29–31]. Sharan et al. [30] collapsed dynamic networks to predict temporal links by summing the matrices associated with networks, thereby saving running time by sacrificing accuracy. To fully utilize topology structure, the Katz index predicts temporal links by counting the number of paths. Matrix decomposition-based algorithms [29,31], such as singular value decomposition (SVD) and tensor decomposition (TD), have been developed to predict temporal links using low-rank approximation. These algorithms initially collapse temporal networks and then predict temporal links based on the collapsed network, which eliminates critical information hidden in dynamic networks, thereby affecting the performance of algorithms. To avoid the collapse of temporal networks, Acar et al. [29] provided a TD method for the temporal link prediction problem, and this method dramatically improved the accuracy of algorithms.

Although considerable efforts have been devoted to the temporal link prediction, some problems remain unsolved, including determining the theoretical relationship among matrix decomposition algorithms and improving the accuracy of algorithms. In this work, the first problem is addressed by proving the equivalence between the eigendecomposition (ED) and nonnegative matrix factorization (NMF) based algorithms. To address the second issue, an NMF-based algorithm is developed without collapsing dynamic networks, which significantly improves the accuracy of algorithms.

Overall, the main contributions of this study can be summarized as follows.

- We prove the equivalence between the eigendecomposition and nonnegative matrix factorization algorithms in temporal networks, which serves as the theoretical foundation for designing NMF-based algorithms for the temporal link prediction problem.
- Two NMF-based frameworks for the temporal link prediction problem have been proposed based on the proven equivalence by using graph communicability. The two frameworks differ greatly in terms of objects to collapse. The first framework collapses temporal features, whereas the second framework collapses temporal networks.
- The proposed method outperforms state-of-the-art methods by using both the artificial and real-world dynamic networks.

The remainder of this paper is organized as follows. The preliminaries are presented in Section 2. Related works are reviewed in Section 3. The equivalence relationship is proven in Section 4. The proposed algorithm is described in Section 5. The experimental results are presented in Section 6. The extension of algorithms and conclusion are provided in Sections 7 and 8, respectively.

## 2. Preliminaries

Terminologies that are extensively used in the subsequent sections are first introduced prior to presenting the detailed description of the proposed algorithms. Let  $\{1, 2, \dots, T\}$  be a finite set of time points. For a given variable, the attached subscript  $t$  represents the value of the variable at time point  $t$  (time  $t$  for short). The *temporal (dynamic) network*  $\mathcal{G}$  is defined as a sequence of networks  $\mathcal{G} = \{G_1, G_2, \dots, G_T\}$ , where  $G_t$  is the network at time  $t$  with a vertex set  $V_t$  and an edge set  $E_t$ . Without loss of generality, we assume that all of the networks in  $\mathcal{G}$  have the same vertex set, i.e.  $G_t = (V, E_t)$ . The temporal network  $\mathcal{G}$  can be represented by a 3-dimensional matrix (tensor)  $W = (w_{ijt})_{n \times n \times T}$ , where  $n$  is the number of vertices (i.e.  $n = |V|$ ) and  $w_{ijt}$  is the weight on edge  $(v_i, v_j)$  in  $G_t$ . Actually,  $W = [W_1, W_2, \dots, W_T]$ , where  $W_t = (w_{ijt})_{n \times n}$  is the weighted adjacency matrix of network  $G_t$  (also called the  $t$ th slice of  $W$ ). The degree of vertex  $v_i$  in network  $G_t$  is defined as the sum of the weights on the edges connected to it, i.e.  $d_{it} = \sum_j w_{ijt}$ . We assume all of the network in  $\mathcal{G}$  are undirected. The *characteristic polynomial* of matrix  $W_t$  is defined as  $P_{G_t}(x) = \det(xI - W_t)$ . The eigenvalue  $\lambda$  satisfies equation  $W_t \mathbf{x} = \lambda \mathbf{x}$  for certain non-zero vectors  $\mathbf{x} \in R^n$ , where  $\mathbf{x}$  is called an eigenvector of matrix  $W_t$  belonging to eigenvalue  $\lambda$ . We denote the eigenvalues of  $G_t$  as  $\lambda_{1t}, \dots, \lambda_{nt}$ . Because  $W_t$  is real and symmetric, the eigenvalues are real numbers. Without loss of generality, we assume that  $\lambda_{1t} \geq \dots \geq \lambda_{nt}$ .

The temporal link prediction problem is an extension of missing link prediction, which is defined as follows: given the temporal network  $\mathcal{G} = \{G_1, G_2, \dots, G_T\}$ , links in network  $G_{T+1}$  are predicted based on  $\mathcal{G}$ , i.e. a function  $f$  should be constructed, such that

$$W_{T+1} = f(W_1, \dots, W_T). \quad (1)$$

## 3. Related works

In this section, we briefly review the matrix-based algorithms for temporal link prediction problem, which are classified into three classes: network collapse, topology, and matrix decomposition-based approaches.

Typical network collapsing-based approaches include collapsing tensor (CT) [32] and weighted CT (WCT) [30]. The CT algorithm collapses  $\mathcal{G}$  by averaging link weights, i.e.

$$X = \sum_{t=1}^T W_t / T. \quad (2)$$

Then, it predicts temporal links by setting  $W_{T+1} = X$ , which is criticized for its assumption that all networks are equally important. In fact, networks close to time  $T + 1$  are more important than those that are far away from it. To overcome this problem, WCT assigns a weight to each network and collapses  $\mathcal{G}$  by damping the time points backward as follows:

$$X = \sum_{t=1}^T (1 - \theta)^{T-t} W_t, \quad (3)$$

where  $\theta \in (0, 1)$  is a parameter controlling the relevant importance of  $W_t$ . WCT predicts temporal links by setting  $W_{T+1} = X$ . Similar strategy is presented [30].

Network collapse-based algorithms are criticized for their low accuracy because they only use the average edge weight, and thus, fail to fully explore the topological structure of temporal networks. To overcome this problem, the Katz index [33] measures the similarity between vertex  $v_i$  and  $v_j$  as the weighted sum of the paths connecting them, i.e.

$$k_{ijt} = \sum_{i=1}^{+\infty} \beta^i p_{ijt}^{(i)}, \quad (4)$$

Download English Version:

<https://daneshyari.com/en/article/4969539>

Download Persian Version:

<https://daneshyari.com/article/4969539>

[Daneshyari.com](https://daneshyari.com)