

一种融合节点变化信息的动态社区发现方法

贺超波¹, 成其伟², 程俊伟^{1*}, 刘星雨¹, 余 鹏¹, 陈启买¹

(1. 华南师范大学计算机学院, 广东广州 510631; 2. 维沃移动通信有限公司, 广东东莞 523859)

摘 要: 动态社区发现旨在检测动态复杂网络中蕴含的社区结构, 对于揭示网络的功能及演化模式具有重要研究价值. 由于相邻时刻网络的社区结构具有平滑性, 前一时刻网络的社区划分信息可以用于监督当前时刻网络的社区划分过程, 但已有方法均难以有效提取这些信息来提高动态社区发现性能. 针对该问题, 提出一种融合节点变化信息的动态社区发现方法 (Semi-supervised Nonnegative Matrix Factorization combining Node Change Information, NCI-SeNMF). NCI-SeNMF 首先采用 k -core 分析方法提取前一时刻社区网络的 degeneracy-core, 并选取 degeneracy-core 中的节点构造社区隶属先验信息, 然后对相邻时刻网络的节点局部拓扑结构变化程度进行量化, 并将其用于进一步修正社区隶属先验信息, 最后通过半监督非负矩阵分解模型集成社区隶属先验信息进行动态社区发现. 在多个人工合成动态网络和真实世界动态网络上进行大量对比实验, 结果表明, NCI-SeNMF 比现有动态社区发现方法在主要评价指标上至少提升了 4.8%.

关键词: 动态社区发现; 半监督非负矩阵分解; k -core 分析; 社区网络; 复杂网络

基金项目: 国家自然科学基金 (No.62077045)

中图分类号: TP311

文献标识码: A

文章编号: 0372-2112(2024)08-2786-13

电子学报 URL: <http://www.ejournal.org.cn>

DOI: 10.12263/DZXB.20221142

A Dynamic Community Discovery Method via Fusing Node Change Information

HE Chao-bo¹, CHENG Qi-wei², CHENG Jun-wei^{1*}, LIU Xing-yu¹, YU Peng¹, CHEN Qi-mai¹

(1. School of Computer Science, South China Normal University, Guangzhou, Guangdong 510631, China;

2. Vivo Mobile Communication Co., Ltd., Dongguan, Guangdong 523859, China)

Abstract: Dynamic community discovery aims to detect community structure in dynamic complex networks, and has important research value for revealing the functions and evolution patterns of networks. Because the community structure of the adjacent snapshot networks is smooth, the community discovery result of the previous snapshot network can be used to supervise the community discovery process of the current snapshot network. However, existing methods are difficult to effectively extract these information to improve the performance of dynamic community discovery. In view of this, a dynamic community discovery method named NCI-SeNMF (Semi-supervised Nonnegative Matrix Factorization combining Node Change Information) is proposed, which can fuse node change information. NCI-SeNMF firstly uses k -core analysis method to extract the degeneracy core of every community network at the previous snapshot, and selects the nodes in the degeneracy core to construct the prior community membership information. Then, it quantifies the change degree of the local topology structure of the nodes in the adjacent snapshot networks, and applies it to further improve the prior community membership information. Finally, it integrates the prior community membership information through semi-supervised nonnegative matrix factorization model to discover dynamic communities. Extensive comparative experiments have been conducted on several synthetic and real-world dynamic networks, and the results show that NCI-SeNMF improves at least 4.8% in term of core evaluation metrics comparing with the existing dynamic community discovery methods.

Key words: dynamic community discovery; semi-supervised nonnegative matrix factorization; k -core analysis; community network; complex networks

Foundation Item(s): National Natural Science Foundation of China (No.62077045)

1 引言

现实世界的复杂网络(例如社交网络、合著关系网络及蛋白质交互网络等)不仅具有隐式的统计特性(例如小世界和无标度^[1]),而且还具有可直接观察的显式特性:社区结构. 社区是指复杂网络中的节点集簇,表现为相同社区的节点间链接紧密,不同社区的节点间链接稀疏^[2,3]. 对复杂网络进行社区发现不仅有利于理解网络的功能和结构特征,而且还具有重要的应用价值. 例如,社交网络的社区被解释为具有相同兴趣爱好的用户群体,可以用于精准的产品推荐^[4]. 论文合著关系网络的社区被理解为具有相同研究兴趣的科研团队,可以为促进更广泛的科研合作推荐团队或研究人员^[5]. 蛋白质交互网络的社区被定义为蛋白质功能模块,可以用于辅助疾病诊断与治疗^[6].

作为一个涉及社会科学、物理学、计算机科学等多个学科的研究话题,社区发现已吸引了大量关注,各类社区发现方法已被大量提出,其中包括基于模块度优化的方法^[7]、基于概率生成模型的方法^[8]、基于标签传播的方法^[9]、基于博弈论的方法^[10]、基于非负矩阵分解的方法^[11]以及基于深度学习的方法^[12]等. 需要指出的是,大部分现有社区发现方法都只能适用于静态复杂网络. 然而,现实中的复杂网络通常是动态变化的,节点与边会随着时间的推移而生成或消失,潜在的社区结构也会随之不断演化^[13]. 适用于动态网络的社区发现(以下简称动态社区发现)方法不仅能发现网络在各个时刻的社区结构,还能进一步分析社区演化模式和驱动因素,最终可进一步拓展社区发现的实际应用范围(如预测网络的演化趋势^[14]). 为此,研究有效的动态社区发现方法具有更迫切的现实需求和更高的应用价值.

目前,动态社区发现已成为复杂网络分析领域的研究热点,并提出了一些代表性方法,这些方法通常可以划分为三类:基于两阶段的方法、基于增量学习的方法以及基于演化聚类的方法. 基于两阶段的方法(如 CPM(Critical Path Method)^[15]和 OCDA(Open Compound Domain Adaptation)^[16])核心思想是第一个阶段先对每一个时刻的网络采用静态社区发现算法独立地进行社区划分,第二个阶段再在不同的时刻上对社区进行匹配,从而可以捕捉社区动态变化过程. 基于增量学习的方法可以在初始时刻网络应用聚类算法检测社区,后续时刻网络只需要分析增量部分对社区结果进行动态更新. 这类方法具有效率优势,如代表性方法 Lime (Low-cost and incremental learning for dynamic heterogeneous information networks)^[17]、DISCAN (Distributed and Incremental Structural Clustering Algorithm for Networks)^[18]和 DyPerm (Permanence for Dynamic commu-

nity detection)^[19]等都具有近似线性的时间复杂度. 基于演化聚类的方法采用双目标优化框架,第一个优化目标为最大化当前时刻的社区发现准确性,第二个优化目标为最小化前后时刻的社区划分差异. 由于能够灵活集成更多有用的约束信息,这类方法往往具有更好的社区发现准确性,同时易采用各种聚类优化模型进行扩展实现,如基于粒子群优化的 DYN-MODPSO (Multi-Objective Discrete Particle Swarm Optimization for DYNamic network)^[20]、基于非负矩阵分解的 DGR-NMF (Dynamic Graph Regularized symmetric Nonnegative Matrix Factorization)^[21]和 Cr-ENMF (Co-regularized Non-negative Matrix Factorization)^[22]等. 虽然已有的动态社区发现方法都具有不同程度的有效性,但它们未能有效利用隐含的先验信息进一步提升性能. 事实上,社区在动态演化的过程中具有平滑性,相邻时刻的网络存在稳定的节点社区隶属信息. 例如,论文合著关系网络中的科研团队在演化过程中通常具有较为稳定的核心成员. 显然,这些信息可以很自然地作为先验信息用于监督当前时刻网络的社区发现过程. 虽然已有方法尝试集成这些社区隶属先验信息来提高性能,如 DGR-NMF (Dynamic Graph Regularized symmetric Nonnegative Matrix Factorization)^[21]和 Cr-ENMF^[22],但如何有效地提取并集成利用仍然需要进一步研究. 具体而言,目前大部分已有方法都无法有效地解决以下两个问题:

(1)由于无论使用何种方法都很难保证社区发现结果是完全正确的,因此从前一时刻网络的社区划分结果提取先验信息时,很容易包含过多错误的社区划分信息.

(2)某些节点随着网络的演化,可能会发生社区转移. 因而即使这些节点在前一时刻网络的社区划分正确,但如果其社区隶属关系在当前时刻发生了转移,与这些节点相关的社区隶属先验信息将是无效的.

显然,错误和无效的社区隶属先验信息将会极大地误导当前时刻网络的社区发现过程并降低性能. 基于此,本文提出了一种融合节点变化信息的半监督非负矩阵分解动态社区发现方法 NCI-SeNMF (Semi-supervised Nonnegative Matrix Factorization combining node Change Information). 为有效处理社区隶属先验信息,NCI-SeNMF 首先引入 k -core 分析方法提取前一时刻各个社区网络的 degeneracy-core. 接着使用 degeneracy-core 的节点来构造先验信息,从而过滤潜在的错误社区划分节点. 然后对节点相邻时刻网络的局部结构变化程度进行量化,将其用于修正已发生社区转移节点的先验信息. 最后使用演化聚类框架设计 NCI-SeNMF 动态社区发现模型,通过图正则非负矩阵分解集成利用社区隶属先验信息以实现动态社区发现性能的提升.

2 社区隶属先验信息处理

本节首先介绍 k -core 分析^[23]相关概念的定义,然后基于 k -core 对社区隶属节点进行特征分析,最后提出前一时网络社区隶属先验信息的提取及修正方法.

2.1 基于 k -core 的社区隶属节点特征实证分析

定义 1 给定一个图 $G=(V,E)$, 一个非负整数 k , $S=(V',E')$ 为 G 的子图, 若 S 同时满足以下两个约束, 则称 S 为 G 的 k -core:

约束 1 对于任意的 $v \in V'$, 都有 $\text{degree}(v) \geq k$, 其中 $\text{degree}(v)$ 表示节点 v 的度.

约束 2 S 是满足约束 1 的最大子图, 即不存在任何 S 的超图 S' 满足约束 1.

定义 2 给定一个图 $G=(V,E)$, 节点 $v \in V$, 一个非负整数 i , 若节点 v 存在于图 G 的 i -core, 但不存在于图 G 的 $(i+1)$ -core, 则称 i 为节点 v 的 coreness.

定义 3 给定一个图 $G=(V,E)$, 一个非负整数 k , 若 G 的 k -core 存在, 但 G 的 $(k+1)$ -core 不存在, 则称 k -core 为 G 的 degeneracy-core.

基于 k -core 相关概念, 采用人工合成网络生成工具 LFR^[24] 生成一个网络 G , 并对其社区隶属节点的特征进行实证分析. G 的节点个数为 500、社区个数为 6、混合参数 μ 为 0.4. 首先使用常用的 Louvain 算法^[25] 对 G 进行社区检测, 最终检测的 6 个社区分别记作 C_1, C_2, \dots, C_6 . 特别地, 将每一个社区 C_x 包含的节点和其节点之间的边看作一个独立的网络, 并表示为社区网络 G_{C_x} . 由于 LFR 生成的网络包含所有节点的真实社区标签, 所以根据是否被划分到正确的社区, 社区网络中的节点可以被分为两类: 正确节点与错误节点. 对每一个社区网络的正确节点集合与错误节点集合分别进行度分布统计, 其中社区网络 G_{C_x} 正确节点集合的度分布定义为

$$P_{\text{correct}}(d, C_x) = \frac{Q_{\text{correct}}(d, C_x)}{N_{\text{correct}}(C_x)} \quad (1)$$

其中, $Q_{\text{correct}}(d, C_x)$ 表示社区网络 G_{C_x} 正确节点集合中度为 d 的节点个数, $N_{\text{correct}}(C_x)$ 表示社区网络 G_{C_x} 正确节点集合的节点总数. 错误节点集合的度分布可以按类似的方法计算.

最终 G 中各个社区网络的度分布计算结果如图 1(a) 所示. 通过图 1(a) 可以直观地看出: 对于错误节点集合, 节点的度大部分都为 1~3, 极少部分度大于等于 4; 而对于正确节点集合, 度大于等于 4 的节点则占了该集合的绝大部分. 可以认为: 错误节点大概率为小度数节点, 而大度数节点几乎不可能为错误节点. 需要说明的是, 在其他网络上发现了该现象, 为避免重复, 这里忽略了相类似的统计结果. 本文认为产生该现象的原因是: 节点在社区网络中的度越大, 则该节点与同一社区其他节点存在的链接就越多. 进而大度数节点相互构成的网络具备高稠密性, 符合社区内部链接紧密的特征, 因此大度数节点的划分结果会趋于正确. 而小度数节点与同一社区其他节点的连接十分稀疏, 其组成的网络不能较好地满足社区链接紧密的特征, 所以错误节点几乎都是小度数节点.

文献[21]指出节点的 coreness 与节点的度数呈强正相关, 并且节点的度数为节点 coreness 的上界, 因此容易推测错误节点具备小 coreness, 而正确节点具备大 coreness. 为验证该推测, 根据式(1), 对每一个社区网络的正确节点与错误节点进行节点 coreness 的分布统计, 此时 $Q_{\text{correct}}(d, C_x)$ 则是表示社区网络 G_{C_x} 的正确节点集合中 coreness 为 d 的节点个数. 最终得到的结果如图 1(b) 所示, 从该图中可以看出, 两类节点的 coreness 分布相对于度分布变得更为清晰. 其中所有社区网络正确节点的 coreness 大部分都为 4, 同时 4 是所有社区网络节点的最大 coreness 值, 即正确节点大部分具有所处社区网络的最大 coreness 值. 相比之下, 大部分错误节点的 coreness 值都是小于等于 3.

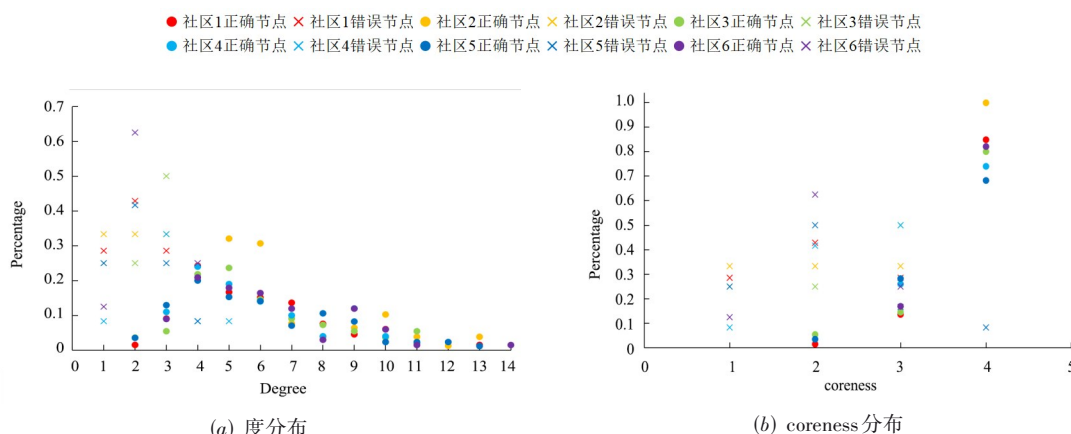


图1 正确和错误节点的度分布与 coreness 分布

综上分析并结合定义 2 与定义 3, 可以发现: degeneracy-core 的节点都具有最大的 coreness, 所以可以认为 degeneracy-core 中的节点基本都为正确节点. 为进一步验证该结论的正确性, 对 G 的 6 个社区网络从 k -core 分层的角度进行节点社区划分准确率统计, 结果如表 1 所示.

表 1 各个社区网络的 k -core 节点社区划分准确率 单位: %

| 社区网络 | 1-core | 2-core | 3-core | 4-core |
|--------|--------|--------|--------|--------|
| 社区网络 1 | 90.41 | 92.96 | 97.01 | 100 |
| 社区网络 2 | 92.86 | 95.12 | 97.50 | 100 |
| 社区网络 3 | 93.22 | 94.83 | 96.30 | 100 |
| 社区网络 4 | 89.29 | 90.09 | 94.34 | 100 |
| 社区网络 5 | 87.63 | 90.43 | 96.47 | 98.31 |
| 社区网络 6 | 89.33 | 90.54 | 97.10 | 100 |

通过表 1 可以看出, 随着 k -core 的 k 值增大, 所包含节点的社区划分准确率随之提升, 并且 5 个社区网络的 degeneracy-core (即 4-core) 的节点划分准确率提升到了 100%, 另外 1 个社区的准确率也提升到了 98.31%. 由此可知, 通过 k -core 分解提取社区网络的 degeneracy-core 可以有效地过滤错误节点, 并尽可能地保留正确节点. 进一步地, 文献[26]指出, 节点的 coreness 越大, 该节点离开网络的概率越小, 因而还可以认为 degeneracy-core 的节点在网络动态变化的过程中离开社区的概率很小, 即 degeneracy-core 的节点同时还具有比较稳定的社区隶属关系.

2.2 先验信息提取

Must-link 约束常被用作先验信息, 用于监督社区发现过程以提高社区划分性能. 若两个节点属于 must-link 约束, 则认为这两个节点隶属于同一个社区. 根据前一节的分析可知, degeneracy-core 的节点基本都具有正确且稳定的社区隶属关系, 因此可以提取前一个时刻各个社区网络 degeneracy-core 的节点社区隶属关系构建 must-link 约束, 并用于监督当前时刻网络的社区发现过程. 具体而言, 如果在上一个时刻两个节点同时在同一个社区网络的 degeneracy-core 中, 则认为这两个节点在当前时刻也会隶属于同一个社区, 从而它们可以构成一个 must-link 约束. 所有的 must-link 约束可以表示为一个约束矩阵 M_t , 其中分量定义如下:

$$M_{ij,t} = \begin{cases} 1, & \text{if } v_i \in \text{degeneracy}(G_{c_x,t-1}) \text{ and } \\ & v_j \in \text{degeneracy}(G_{c_x,t-1}) \\ 0, & \text{otherwise} \end{cases} \quad (2)$$

其中, $\text{degeneracy}(G_{c_x,t-1})$ 表示第 $(t-1)$ 时刻社区网络 G_{c_x} 的 degeneracy-core.

处于相同社区的节点应具有更为相似的社区隶属关系表示, 为此对于任意节点 v_i 和 v_j , 它们的 must-link 约束应能限制它们在 t 时刻对应的社区隶属关系表示向量 $H_{i,t}$ 和 $H_{j,t}$ 之间的距离, 即有

$$\begin{aligned} R(H_t) &= \sum_{i=1}^N \sum_{j=1}^N \|M_{ij,t} H_{i,t} - H_{j,t}\|_F^2 \\ &= \sum_{i=1}^N \sum_{j=1}^N M_{ij,t} (H_{i,t} H_{i,t}^T + H_{j,t} H_{j,t}^T - 2 H_{i,t} H_{j,t}^T) \\ &= 2 \sum_{i=1}^N \mathcal{D}(M_t) H_{i,t} H_{i,t}^T - 2 \sum_{i=1}^N \sum_{j=1}^N M_{ij,t} H_{i,t} H_{j,t}^T \\ &= 2 \text{tr}(H_t^T \mathcal{D}(M_t) H_t) - 2 \text{tr}(H_t^T M_t H_t) \end{aligned} \quad (3)$$

其中, N 为节点数量; $\text{tr}(\cdot)$ 表示矩阵的迹; 函数 $\mathcal{D}(M_t)$ 以矩阵 M_t 作为输入, 输出一个对角矩阵:

$$\mathcal{D}(M_t)_{ij} = \begin{cases} \sum_{j=1}^N M_{ij,t}, & \text{if } i=j \\ 0, & \text{if } i \neq j \end{cases} \quad (4)$$

2.3 先验信息修正

尽管 M_t 包含高度可靠的先验信息, 但仍可能存在由于节点发生转移而导致相关先验信息失效的问题. 即对于前后两个相邻时刻, 若某个节点发生了社区转移, 即使该节点相关的先验信息在上一时刻是正确的, 但却不适用于当前时刻. 以图 2 所示的网络为例, 在时刻 t 中, 编号 1 至 5 的节点属于 1 号社区, 编号 6 至 10 的节点属于 2 号社区. 在时刻 $(t+1)$ 中, 2 号节点从原先的 1 号社区转移到了 2 号社区中. 假设提供了符合上一时刻的先验信息, 即 2 号与 1、3、4、5 号节点同属于 1 号社区, 那么可能会导致在时刻 $(t+1)$ 由于先验信息的监督作用错误地将 2 号节点划分到社区 1 中. 基于此, 有必要减弱发生社区转移的节点对当前时刻社区划分所产生的消极影响.

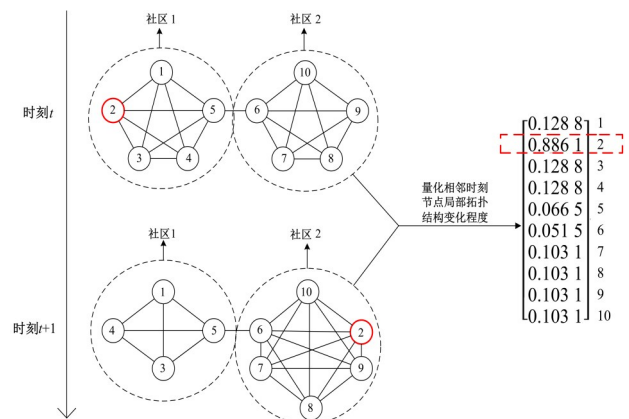


图 2 相邻时刻节点发生社区转移的网络

直观地,若一个节点发生了社区转移,则当前时刻以该节点为中心的局部网络结构相比于前一个时刻会产生剧烈的变化,而变化程度可以使用下式进行计算:

$$\text{change}(i, t) = 1 - \text{similarity}(\mathbf{s}_{i,t-1}, \mathbf{s}_{i,t}) \quad (5)$$

其中, $\mathbf{s}_{i,t}$ 是一个关于节点 i 的相似度向量, 其元素 $s_{ij,t}$ 表示在时刻 t 中节点 i 与节点 j 的相似度. $\text{similarity}(\cdot)$ 则是度量两个向量相似度的函数.

本文选择 RA 相似度^[27]来构造节点的局部相似性向量:

$$s_{ij,t} = \sum_{v \in \Gamma(i) \cap v \in \Gamma(j)} \frac{1}{\text{degree}(v)} \quad (6)$$

其中, $\Gamma(i)$ 表示节点 i 的直接邻居, $\text{degree}(v)$ 表示节点 v 的度. 对于 similarity 函数则选择余弦相似度:

$$\text{similarity}(\mathbf{s}_{i,t-1}, \mathbf{s}_{i,t}) = \frac{\sum_{k=1}^N (s_{ik,t-1} \times s_{ik,t})}{\sqrt{\sum_{k=1}^N (s_{ik,t-1})^2} \times \sqrt{\sum_{k=1}^N (s_{ik,t})^2}} \quad (7)$$

通过式(5), 可得到相邻时刻节点局部网络结构的变化程度. 若变化程度越大, 则表明该节点有很大概率发生了社区转移. 若变化程度较小, 则表明该节点发生社区转移的概率很小. 以图2左边的网络为例子, 使用式(5)计算各个节点的局部结构变化程度向量, 得到的结果如图2右边的向量所示, 可以看出节点2的局部结构变化程度为0.8861, 远大于其他节点, 说明该节点大概率已经发生了社区转移. 将相邻时刻节点局部网络结构的变化程度用于修正先验信息, 若节点局部结构变化程度越强, 则说明发生社区转移的可能性就越大, 因此应尽量削弱与该节点相关的先验信息, 反之若变化程度越弱, 则说明该节点发生社区转移的可能性就越小, 应尽量保留与该节点相关的先验信息. 为此, 结合式(5)将式(3)重写为

$$\begin{aligned} R(\mathbf{H}_t) &= \sum_{i=1}^N \sum_{j=1}^N M_{ij,t} S_{ii,t} S_{jj,t} (\mathbf{H}_{i,t} \mathbf{H}_{i,t}^T \\ &\quad + \mathbf{H}_{j,t} \mathbf{H}_{j,t}^T - 2\mathbf{H}_{i,t} \mathbf{H}_{j,t}^T) \\ &= 2 \sum_{i=1}^N \mathcal{D}(\mathbf{S}_t \mathbf{M}_t \mathbf{S}_t)_{ii} \mathbf{H}_{i,t} \mathbf{H}_{i,t}^T \\ &\quad - 2 \sum_{i=1}^N \sum_{j=1}^N M_{ij,t} S_{ii,t} S_{jj,t} \mathbf{H}_{i,t} \mathbf{H}_{j,t}^T \\ &= 2\text{tr}(\mathbf{H}_t^T \mathcal{D}(\mathbf{S}_t \mathbf{M}_t \mathbf{S}_t) \mathbf{H}_t) - 2\text{tr}(\mathbf{H}_t^T \mathbf{S}_t \mathbf{M}_t \mathbf{S}_t \mathbf{H}_t) \\ &= 2\text{tr}(\mathbf{H}_t^T (\mathcal{D}(\mathbf{S}_t \mathbf{M}_t \mathbf{S}_t) - \mathbf{S}_t \mathbf{M}_t \mathbf{S}_t) \mathbf{H}_t) \end{aligned} \quad (8)$$

其中, \mathbf{S}_t 是一个对角矩阵, 定义如下:

$$S_{ij,t} = \begin{cases} \text{similarity}(\mathbf{s}_{i,t-1}, \mathbf{s}_{i,t}), & \text{if } i=j \\ 0, & \text{if } i \neq j \end{cases} \quad (9)$$

3 动态社区发现方法 NCI-SeNMF

本节首先介绍 NCI-SeNMF 基于演化聚类框架的动态社区发现模型, 然后推导模型的迭代优化规则, 最后设计相应的动态社区发现算法, 并进行时间复杂度分析.

3.1 动态社区发现模型

基于演化聚类的动态社区发现框架形式上可表示为

$$\text{Cost} = \alpha \text{CS} + (1 - \alpha) \text{CT} \quad (10)$$

其中, CS 是当前时刻网络的社区划分损失项, 用于衡量社区划分结果拟合当前时刻网络结构的好坏程度, 而 CT 是网络的时变损失项, 用于衡量当前时刻的社区结构与前一个时刻社区结构的相似程度, 参数 α 则用于控制 CS 和 CT 两个损失项的贡献. 基于演化聚类的动态社区发现框架本质上是一种半监督学习框架, CT 可以作为监督项用于集成来自于前一个时刻网络的先验信息.

基于演化聚类框架, 考虑到对称非负矩阵分解 SNMF^[11,28]具有良好的网络数据聚类能力且易扩展, 因此本文选择 SNMF 构建损失项 CS, 具体步骤:

(1) 给定 T 个时刻的网络快照序列 $G = \{G_1, G_2, \dots, G_T\}$, 令 \mathbf{A}_t 表示第 t 时刻网络 G_t 的邻接矩阵, 对于 $\forall A_{ij,t} \in \mathbf{A}_t$, 如果在第 t 时刻节点 i 与节点 j 之间存在边, $A_{ij,t} = 1$, 否则 $A_{ij,t} = 0$.

(2) 应用 SNMF 将 \mathbf{A}_t 分解成低秩矩阵因子 \mathbf{H}_t 与其转置矩阵的乘积, 即 $\mathbf{A}_t \approx \mathbf{H}_t \mathbf{H}_t^T$, 分解损失表示为

$$\min L(\mathbf{H}_t) = \|\mathbf{A}_t - \mathbf{H}_t \mathbf{H}_t^T\|_F^2 \quad (11)$$

其中, $\|\cdot\|_F$ 为 Frobenius 范数; \mathbf{H}_t 可以作为社区隶属关系表示矩阵, 对于 $\forall H_{ij,t} \in \mathbf{H}_t$, 其表示在第 t 时刻节点 i 隶属于社区 j 的强度.

(3) 使用 $L(\mathbf{H}_t)$ 作为当前时刻网络的损失项, 即 $\text{CS} = L(\mathbf{H}_t)$.

对于时变损失项 CT, 由于其可以视为集成先验信息的监督项, 因此可以很自然地把 $R(\mathbf{H}_t)$ 视为时变损失项 CT. 结合 $L(\mathbf{H}_t)$ 和 $R(\mathbf{H}_t)$, 最终可得到融合节点变化信息的半监督非负矩阵分解动态社区发现模型:

$$\begin{aligned} \min \mathcal{O}(\mathbf{H}_t) &= \alpha \|\mathbf{A}_t - \mathbf{H}_t \mathbf{H}_t^T\|_F^2 + (1 - \alpha) \text{tr}(\mathbf{H}_t^T (\mathcal{D}(\mathbf{S}_t \mathbf{M}_t \mathbf{S}_t) \\ &\quad - \mathbf{S}_t \mathbf{M}_t \mathbf{S}_t) \mathbf{H}_t) \\ \text{s.t. } H_{ij,t} &\geq 0, \forall i, j \end{aligned} \quad (12)$$

需要注意的是, 当 $\alpha = 1$ 或 $t = 1$ 时, 由于缺乏先验约束, 可以直接使用 SNMF 进行社区发现.

3.2 优化求解

通过最小化目标函数 $\mathcal{O}(\mathbf{H}_t)$ 可获得 \mathbf{H}_t 的近似解, 从而获得社区划分结果. 由于 \mathbf{H}_t 带有非负值约束, 最小化 $\mathcal{O}(\mathbf{H}_t)$ 可以转化为受限约束求极值问题, 这可以应用拉格朗日乘数方法进行优化求解. 设 Ψ 为约束 $\mathbf{H}_t \geq 0$ 的拉格朗日乘子, 式(12)的拉格朗日函数可表示为

$$L(\mathbf{H}_t) = \alpha \|\mathbf{A}_t - \mathbf{H}_t \mathbf{H}_t^T\|_F^2 + (1-\alpha) \text{tr}(\mathbf{H}_t^T (\mathbf{D}(\mathbf{Q}_t) - \mathbf{Q}_t) \mathbf{H}_t) + \text{tr}(\Psi \mathbf{H}_t) \quad (13)$$

其中, $\mathbf{Q}_t = \mathbf{S}_t \mathbf{M}_t \mathbf{S}_t^T$.

根据矩阵 Trace 与 Frobenius 范数的关系: $\|\mathbf{X}\|_F^2 = \text{tr}(\mathbf{X}\mathbf{X}^T)$, 式(13)可以重写为

$$L(\mathbf{H}_t) = \alpha (\text{tr}(\mathbf{A}_t \mathbf{A}_t^T) - 2\text{tr}(\mathbf{A}_t \mathbf{A}_t^T) + \text{tr}(\mathbf{H}_t \mathbf{H}_t^T \mathbf{H}_t \mathbf{H}_t^T)) + (1-\alpha) \text{tr}(\mathbf{H}_t^T (\mathbf{D}(\mathbf{Q}_t) - \mathbf{Q}_t) \mathbf{H}_t) + \text{tr}(\Psi \mathbf{H}_t^T) \quad (14)$$

$L(\mathbf{H}_t)$ 对 \mathbf{H}_t 求偏导可得到下式:

$$\frac{\partial L}{\partial \mathbf{H}_t} = -4\alpha \mathbf{A}_t \mathbf{H}_t - 2(1-\alpha) \mathbf{Q}_t \mathbf{H}_t + 4\alpha \mathbf{H}_t \mathbf{H}_t^T \mathbf{H}_t + 2(1-\alpha) \mathbf{D}(\mathbf{Q}_t) \mathbf{H}_t + \Psi \quad (15)$$

根据 KKT 条件 $\Psi_{ij} H_{ij,t} = 0$, 可得:

$$\left[-2\alpha (\mathbf{A}_t \mathbf{H}_t)_{ij} - (1-\alpha) (\mathbf{Q}_t \mathbf{H}_t)_{ij} + 2\alpha (\mathbf{H}_t \mathbf{H}_t^T \mathbf{H}_t)_{ij} + (1-\alpha) (\mathbf{D}(\mathbf{Q}_t) \mathbf{H}_t)_{ij} \right] H_{ij,t} = 0 \quad (16)$$

由上式可得 $H_{ij,t}$ 的乘性迭代更新规则为

$$H_{ij,t} = H_{ij,t} \frac{(2\alpha \mathbf{A}_t \mathbf{H}_t + (1-\alpha) (\mathbf{Q}_t \mathbf{H}_t))_{ij}}{(2\alpha \mathbf{H}_t \mathbf{H}_t^T \mathbf{H}_t + (1-\alpha) (\mathbf{D}(\mathbf{Q}_t) \mathbf{H}_t))_{ij}} \quad (17)$$

注意, 当 $T=1$ 时, 由上述过程同理可得 \mathbf{H}_1 的乘性迭代更新规则:

$$H_{ij,1} = H_{ij,1} \frac{(\mathbf{A}_1 \mathbf{H}_1)_{ij}}{(\mathbf{H}_1 \mathbf{H}_1^T \mathbf{H}_1)_{ij}} \quad (18)$$

得到社区隶属强度矩阵 \mathbf{H}_t 后, 可通过 $\arg\max_i H_{i,t}$

获得节点 i 隶属的社区编号.

3.3 社区发现算法及时间复杂度分析

使用 NCI-SeNMF 方法对动态社区发现的流程可通过算法 1 进行描述.

可以分析, 算法 1 的时间消耗主要集中在三个部分. 第一部分 (第 12~17 行) 是提取社区网络的 degeneracy-core, 其时间复杂度为 $O(n+m)$, 其中 n 和 m 分别为网络的节点个数和边数. 第二部分 (第 22

算法 1 NCI-SeNMF 社区发现

输入: T 个时刻的网络序列 $G = \{G_1, G_2, \dots, G_T\}$, 平衡参数 α

输出: T 个时刻的社区划分结果序列 $C = \{C_1, C_2, \dots, C_T\}$

```

1. communities = [] // 保存每个时刻社区划分结果的集合
2. 随机初始化  $\mathbf{H}_1$ 
3. WHILE  $\mathcal{O}(\mathbf{H}_1)$  不收敛 DO
4.   根据式(18)更新  $\mathbf{H}_1$ 
5. END WHILE
6. 根据  $\mathbf{H}_1$  获得时刻 1 的社区划分结果  $C_1$ , 添加到 communities 集合中
7. FOR  $t = 2$  to  $t = T$ 
8.   FOR  $\forall$  community in  $C_{t-1}$ :
9.     根据 community 构造社区网络  $\mathcal{G}$ 
10.     $i \leftarrow 1$  // 该变量表示  $k$ -core 的  $k$  值
11.    degeneracy-core  $\leftarrow$  null
12.    WHILE  $\mathcal{G}$  的节点个数不为 0 DO
13.      WHILE  $\mathcal{G}$  不满足  $k$  值为  $i$  的  $k$ -core DO
14.        删除度数小于  $i$  的节点及该节点关联的边
15.        更新其余节点的度数
16.      END WHILE
17.    END WHILE
18.    degeneracy-core  $\leftarrow \mathcal{G}$ 
19.     $i \leftarrow i + 1$ 
20.    根据式(2)对矩阵  $\mathbf{M}$  赋值
21.  End FOR
22. 根据式(6)、式(7)和式(9)构造矩阵  $\mathbf{S}_t$ 
23. 随机初始化  $\mathbf{H}_t$ 
24.  WHILE  $\mathcal{O}(\mathbf{H}_t)$  不收敛 DO
25.    根据式(17)更新  $\mathbf{H}_t$ 
26.  END WHILE
27.  根据  $\mathbf{H}_t$  获得时刻  $t$  的社区划分结果  $C_t$ , 添加到 communities 集合中
28. END FOR
29. RETURN communities

```

行) 是构造节点相邻时刻的局部结构相似度矩阵 \mathbf{S}_t , 其时间复杂度为 $O(n^2)$. 第三部分是迭代更新社区隶属矩阵 \mathbf{H}_t (第 24~26 行), 其时间复杂度为 $O(t_{\text{iter}} cn^2)$, 其中, c 为平均社区个数, t_{iter} 为平均迭代次数. 那么对于总共有 T 个时刻的网络, 总的时间复杂度为 $O(T(n+m+n^2+t_{\text{iter}} cn^2))$.

4 实验研究

为验证 NCI-SeNMF 的有效性, 本文选取人工合成动态网络和真实动态网络进行实验对比分析. 实验环境为 64 位 Windows 10 操作系统、3.5 GHz Intel Core i9 PC-11900K CPU 和 64 GB RAM, 所有方法都使用 Python 3.7 编程实现.

4.1 评价指标

为评价社区发现方法的性能,对于具有真实社区标签的网络,采用标准化互信息 NMI(Normalized Mutual Information)作为评价指标,对于不具有真实社区标签的网络,采用模块度(Modularity)作为评价指标. NMI 和 Modularity 分别定义如下:

$$\text{NMI}(C, C') = \frac{-2 \sum_{i=1}^q \sum_{j=1}^l F_{ij} \log \left(\frac{F_{ij} N}{F_i \cdot F_j} \right)}{\sum_{i=1}^q F_i \cdot \log \left(\frac{F_i}{N} \right) + \sum_{j=1}^l F_j \cdot \log \left(\frac{F_j}{N} \right)} \quad (19)$$

其中, C 和 C' 分别表示为检测的社区集合和真实的社区集合, F 是一个混淆矩阵, 其元素 F_{ij} 表示同时属于 C_i 和 C'_j 的节点总个数, F_i 表示混淆矩阵 F 的第 i 行元素值的和, F_j 表示混淆矩阵 F 的第 j 列元素值的和, q 和 l 分别表示检测的社区个数和真实的社区个数.

$$\text{Modularity} = \frac{1}{2m} \sum_i \sum_j \left[A_{ij} - \frac{d_i d_j}{2m} \right] \delta_{ij} \quad (20)$$

其中, d_i 表示节点 i 的度数; δ_{ij} 用于判断节点 i 与节点 j 是否属于同一社区, 若属于同一社区, $\delta_{ij} = 1$, 否则 $\delta_{ij} = 0$.

NMI 用于评价社区划分的准确性, 而 Modularity 用于评价社区结构的链接紧密度. 对于 NMI 和 Modularity, 取值越高, 则表明社区划分性能越好.

4.2 对比方法

为了对比验证 NCI-SeNMF 的有效性, 选择下列 5 个具有代表性的动态社区发现方法进行对比.

(1) DNMF (Deep Nonnegative Matrix Factorization)^[29]: DNMF 应用 NMF 检测当前时刻网络的社区, 并通过图正则项约束当前时刻网络的社区划分结果与上一时刻网络的社区划分结果保持相似.

(2) DGR-NMF^[21]: DGR-NMF 基于演化聚类框架, 通过引入网络的几何结构来表示短时间内的网络平滑度, 同时应用图正则化 SNMF 来检测动态社区.

(3) FaceNet (Face recognition and clustering Networks)^[30]: FaceNet 是一种经典的基于生成模型的动态社区发现方法, 它将社区检测和社区演化集成到一个统一的概率模型中, 使得社区演化模式能够约束当前时刻网络的社区检测过程.

(4) ESPRA (Evolutionary clustering based on Structural Perturbation and Resource Allocation similarity)^[31]: ESPRA 是一种结合了结构扰动和拓扑特征的演化聚类方法, 它首先通过融合结构扰动和网络拓扑信息来构造相似度, 然后通过密度聚类来发现社区结构.

(5) Lime^[17]: Lime 基于增量学习框架, 其首先通过共享特征策略学习节点的增量式表示, 然后应用 k -means 聚类算法获得各个时刻网络的社区.

在实验中, NCI-SeNMF 的参数 α 设置为 0.7. 为公平比较, 所有基准方法的超参都设置为最优, 并且每次实验都重复运行 10 次并取平均值作为最终结果.

4.3 人工合成动态网络的对比分析

(1) Dynamic GN 网络

Dynamic GN 是由 Lin 等人^[30]开发的一种基准动态网络生成工具. 使用该工具生成动态网络需要设置 5 个参数, 分别是节点数 n , 社区数 q , 网络序列数 T , 控制社区噪音程度的混合参数 z , 表示节点社区转移概率的参数 nc . z 取值越大, 表明社区结构越模糊. nc 取值越大, 表明相邻时刻的网络结构变化越剧烈. 本文主要通过变更 z 和 nc 两个参数生成 4 个人工合成动态网络, 具体参数的设置如表 2 所示.

表 2 Dynamic GN 网络

| Networks | n | q | T | z | nc |
|----------|-----|-----|-----|-----|------|
| GN1 | 128 | 4 | 10 | 4 | 10% |
| GN2 | 128 | 4 | 10 | 4 | 30% |
| GN3 | 128 | 4 | 10 | 5 | 10% |
| GN4 | 128 | 4 | 10 | 5 | 30% |

将 NCI-SeNMF 与基准方法在 4 个人工合成动态网络进行性能对比, 使用的评价指标为 NMI, 结果如图 3 所示. 可以看出, 在每一个 Dynamic GN 网络, 相对于其他 4 个基准方法, NCI-SeNMF 除了在个别时刻上取得的 NMI 略低于最佳 NMI 外, 大多数时刻取得的 NMI 都要高于其余对比方法的 NMI, 这说明 NCI-SeNMF 在面对社区结构模糊程度不同和变化剧烈程度不同的动态网络都有着更好的社区检测性能. 此外, 在稳定性方面, 可以观察到 NCI-SeNMF 的性能变化曲线是相对平稳的, 而其余 4 个基准方法的性能随着时间的推移都有着不同程度的波动, 说明 NCI-SeNMF 还具有更好的稳定性. 需要指出的是, 在每个网络的时刻 1 中 NCI-SeNMF、DNMF 和 DGR-NMF 具有几乎相同的性能. 这是因为它们在时刻 1 中由于缺乏前一时刻网络的先验信息监督, 都退化为基于 SNMF 的静态社区发现方法, 但它们均优于 FaceNet、ESPRA 及 Lime.

(2) Dynamic LFR 网络

Dynamic LFR 是由 Greene 等人^[32]开发的一种基于 LFR 的动态网络生成工具, 该工具可以生成不同网络演化事件驱动的无向无权网络序列. 本实验选择的网络演化事件包括以下 4 类: 节点隶属社区转移 (Switch)、社区合并与分裂 (MergeSplit)、社区出生与死亡 (Birth-Death) 以及社区扩张与收缩 (ExpandContrat), 在实验中为每类网络演化事件都生成 2 个人工合成动态网络, 最终生成共计 8 个人工合成动态网络, 各个网络的参数设置如表 3 所示.

每一类演化事件的网络生成过程说明如下: 对于

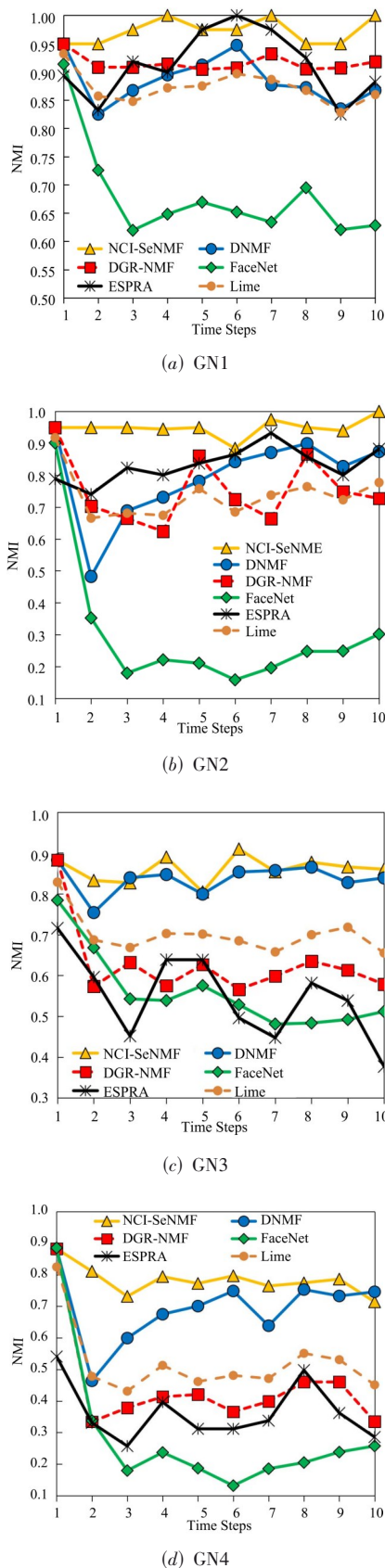


图3 在Dynamic GN网络上的性能比较

表3 Dynamic LFR网络

| Networks | n | t | μ | p | r |
|----------------|-----|-----|-------|-----|-----|
| Switch1 | 500 | 10 | 0.4 | 0.2 | — |
| Switch2 | 500 | 10 | 0.4 | 0.5 | — |
| MergeSplit1 | 500 | 10 | 0.4 | 0.1 | — |
| MergeSplit2 | 500 | 10 | 0.4 | 0.2 | — |
| BirthDeath1 | 500 | 10 | 0.4 | 0.1 | — |
| BirthDeath2 | 500 | 10 | 0.4 | 0.2 | — |
| ExpandContrat1 | 500 | 10 | 0.4 | 0.2 | 0.2 |
| ExpandContrat2 | 500 | 10 | 0.4 | 0.2 | 0.4 |

节点隶属社区转移事件,每个时刻会随机选择 $p \times 100\%$ 的节点更改其隶属社区.对于出生与死亡事件,每个时刻会从其他社区中移除节点来创建为 $p \times 100\%$ 的新社区,并随机删除 $p \times 100\%$ 的社区.对于扩张与收缩事件,每个时刻会随机选择 $p \times 100\%$ 的社区以 $r \times 100\%$ 的大小进行扩张或收缩操作.对于合并与分裂事件,每个时刻会随机选取 $p \times 100\%$ 的社区进行分裂,随机选择 $p \times 100\%$ 的社区进行合并.从表3可以发现,对于基于相同演化事件生成的两个网络,编号为2的网络比编号为1的网络设置有更大的参数 p 或 r ,因此前者比后者发生演化事件的概率更大或者社区变化程度更大,从而导致社区检测难度通常也更高.

图4给出了各方法在8个Dynamic LFR网络的评价结果.首先通过图4(a)~(d)可以看出,在面对发生节点隶属社区转移事件的网络以及发生社区合并与分裂事件的网络时,NCI-SeNMF和DGR-NMF

有着相似的社区发现性能,且两者的社区发现性能在大多数时刻都高于其他3个对比方法.进一步观察图4(e)~(h),可以发现在面对发生社区出生与死亡事件的网络以及发生社区扩张与收缩事件的网络时,DGR-NMF的社区发现性能变得相对较低,而NCI-SeNMF却依旧能保持较好的性能.不仅如此,通过观察图4各方法性能曲线的平稳程度可以发现,NCI-SeNMF不论在具有哪种演化事件的网络,其性能都能够保持相对稳定,而其余方法的性能常常会出现较为剧烈的波动.

为进一步比较各方法在Dynamic LFR网络上的性能稳定性,首先计算各方法在每个网络所有时刻的平均NMI,然后再计算在相同演化事件的两个网络下各方法平均NMI的下降率,最终结果如图5所示.从图5可以看出,当演化事件发生概率增加(p 或 r 增大)时,DGR-NMF、FaceNet、ESPRA及Lime的社区发现性能下降较快,NCI-SeNMF和DNMF则呈现相对较低的下降率.例如,在Switch事件网络中,NCI-SeNMF和DNMF的性能下降率分别为0.47%和0.88%,而DGR-NMF、FaceNet、ESPRA及Lime分别为4.65%、7.66%、4.54%及2.89%.需要指出的是NCI-SeNMF的下降率仅在MergeSplit事件类型网络

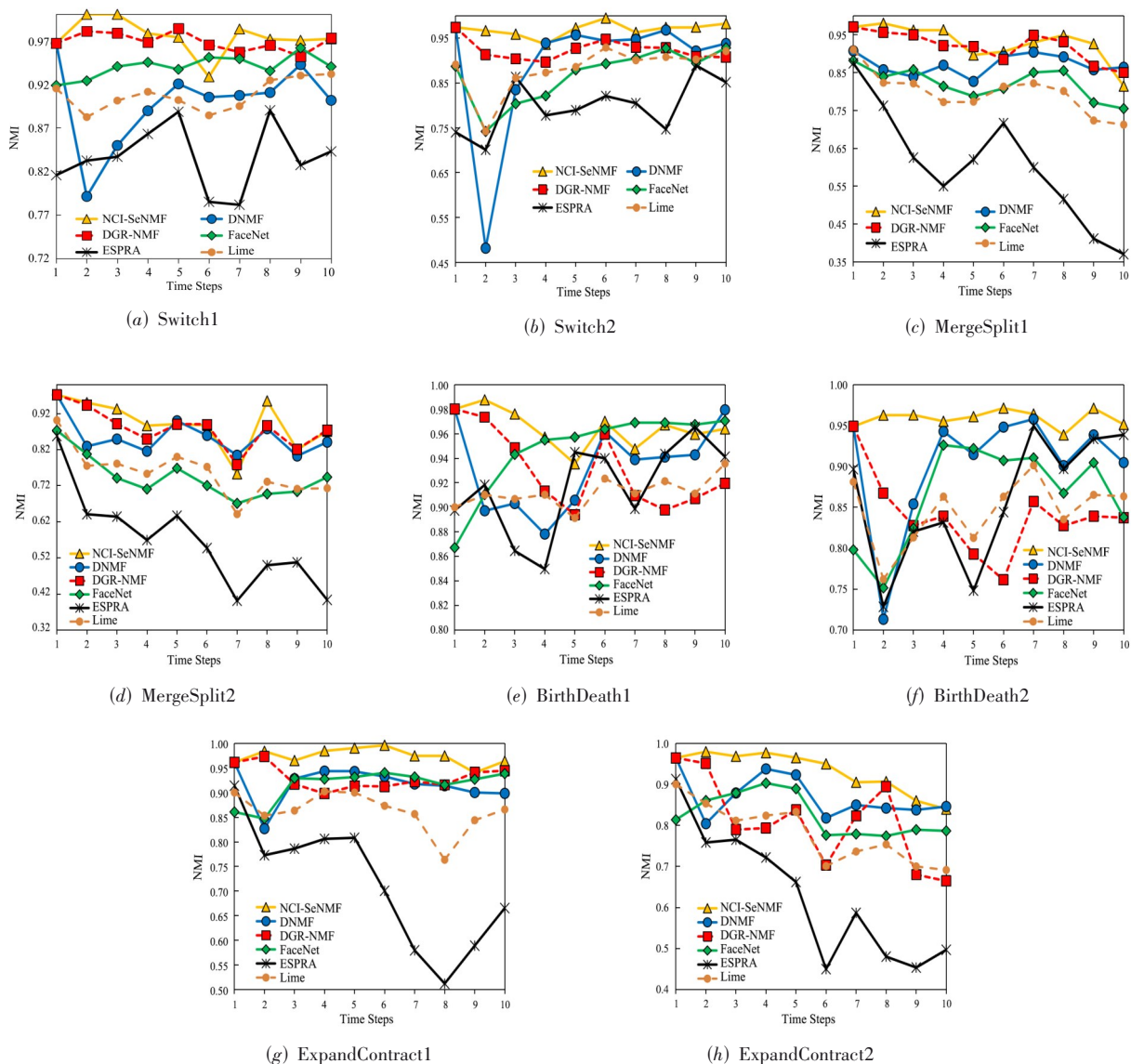


图4 Dynamic LFR 网络上的性能比较

中略高于 DNMF, 说明 NCI-SeNMF 的性能总体更为稳定。

4.4 真实动态网络的对比分析

人工合成动态网络往往只能模拟出某一类型的社区演化事件, 而现实世界的网络往往能同时发生多种类型的社区演化事件. 因此为了进一步验证 NCI-SeNMF 的有效性, 选取了如下 2 个真实世界的动态网络进行实验分析:

(1) Enron 邮件网络^[13]. 该网络来自于美国企业 Enron 的电子邮件集合, 选取了 2001 年 151 名员工的 252 759 封电子邮件数据, 每个月的数据构成一个时刻的网络快照, 共有 12 个时刻的网络快照序列. 网络节点代表员工, 边代表两名员工之间有邮件通信。

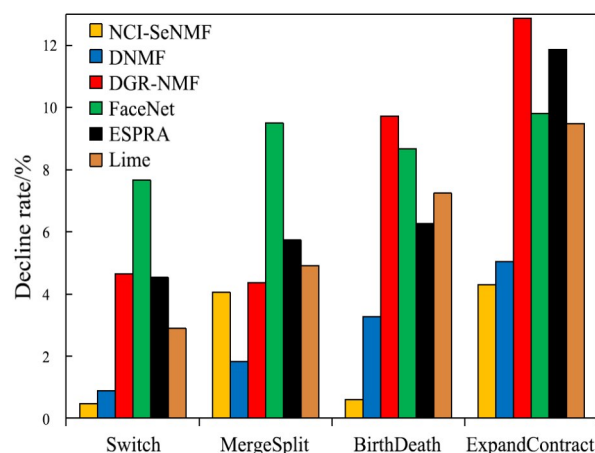


图5 各方法的平均 NMI 下降率

(2)SFHH 会议交互网络^[33]. 该网络描述了 2009 年法国尼斯 SFHH 会议 6 月 4 日至 5 日间 405 名与会者的面对面互动情况,它通过设置相同时间间隔将互动等分为 25 份,然后合并每一份数据,共生成 6 个时刻的网络快照序列.

由于上述网络都不具备真实的社区标签,因此采用的评价指标是模块度. 各方法在 Enron 和 SFHH 动态网络上的性能评价结果分别如表 4 和表 5 所示.

表 4 Enron 网络上的性能比较

| T | NCI-SeNMF | DNMF | DGR-NMF | FaceNet | ESPRA | Lime |
|-----|----------------|----------------|----------------|----------------|----------------|----------------|
| 1 | 0.584 2 | 0.584 2 | 0.584 2 | 0.617 1 | <u>0.600 2</u> | 0.584 2 |
| 2 | 0.645 1 | 0.629 4 | <u>0.635 6</u> | 0.529 0 | 0.624 8 | 0.600 2 |
| 3 | 0.635 7 | 0.588 7 | 0.604 1 | 0.576 2 | <u>0.634 7</u> | 0.591 7 |
| 4 | <u>0.629 1</u> | 0.604 1 | 0.618 8 | 0.581 5 | 0.642 2 | 0.610 5 |
| 5 | 0.542 3 | 0.522 3 | <u>0.526 5</u> | 0.516 9 | 0.471 3 | 0.500 6 |
| 6 | 0.667 3 | 0.612 0 | <u>0.651 2</u> | 0.645 9 | 0.611 0 | 0.602 3 |
| 7 | <u>0.640 1</u> | 0.619 2 | 0.621 9 | 0.591 9 | 0.646 9 | 0.611 3 |
| 8 | 0.516 5 | 0.535 9 | <u>0.527 5</u> | 0.507 2 | 0.507 9 | 0.501 5 |
| 9 | 0.608 6 | <u>0.597 0</u> | 0.601 1 | 0.588 9 | 0.585 3 | 0.576 2 |
| 10 | 0.562 9 | 0.536 5 | 0.544 7 | 0.532 2 | <u>0.552 9</u> | 0.532 9 |
| 11 | <u>0.572 4</u> | 0.561 4 | 0.563 6 | 0.565 2 | 0.581 1 | 0.551 3 |
| 12 | 0.609 7 | 0.573 2 | 0.579 3 | <u>0.583 6</u> | 0.561 0 | 0.573 7 |

注:粗体表示最好结果,下划线表示次优结果.

通过表 4 可以看出,不考虑没有先验信息监督的第一个时刻,NCI-SeNMF 在剩余 11 个时刻中有 7 个时刻都能取得最高的模块度且有 3 个时刻取得次优的模块度,而 DNMF、DGR-NMF、FaceNet、ESPRA 和 Lime 分别有 1 个、0 个、1 个、3 个和 1 个时刻取得最高模块度. 这说明在 Enron 网络中,NCI-SeNMF 相对于另外 4 个对比方法有着更好的性能. 同样地,从表 5 中可以看出,不考虑没有先验信息监督的第一个时刻,NCI-SeNMF 在剩余的 5 个时刻都能取到最高的模块度,并且相对于同一时刻的次优模块度有着较大的提升,分别提升了 12.4%、4.8%、9.8%、15.4% 及 9.2%.

综上分析可知,通过在人工合成动态网络和真实

表 5 SFHH 网络上的性能比较

| T | NCI-SeNMF | DNMF | DGR-NMF | FaceNet | ESPRA | Lime |
|-----|---------------|---------------|---------------|---------------|---------------|--------|
| 1 | <u>0.3066</u> | <u>0.3066</u> | <u>0.3066</u> | 0.2786 | 0.3090 | 0.3012 |
| 2 | 0.3868 | 0.2704 | <u>0.3441</u> | 0.3204 | 0.3031 | 0.2928 |
| 3 | 0.2894 | 0.2690 | <u>0.2761</u> | 0.2099 | 0.1863 | 0.2219 |
| 4 | 0.2527 | <u>0.2302</u> | 0.2052 | 0.2179 | 0.1572 | 0.1973 |
| 5 | 0.3605 | 0.3094 | 0.2759 | <u>0.3124</u> | 0.2815 | 0.2976 |
| 6 | 0.4261 | 0.3846 | <u>0.3903</u> | 0.3842 | 0.3837 | 0.3743 |

注:粗体表示最好结果,下划线表示次优结果.

动态网络进行对比实验,结果表明 NCI-SeNMF 可以有效处理具有各种复杂演化事件的网络,有着更好且稳定的社区发现性能.

4.5 时间效率分析

为测试 NCI-SeNMF 的运行效率,选择 GN1、Switch1、Enron 及 SFHH 网络计算其运行时间,并与其他方法进行对比分析,结果如表 6 所示. 从表 6 中可以看出,基于增量学习的 Lime 方法具有较为明显的运行效率优势,NCI-SeNMF 与其他同样基于 NMF 的方法 DNMF 和 DGR-NMF 相比具有相近时间成本,甚至可以在个别网络(如 Enron)获得次优的运行效率.

表 6 时间效率比较

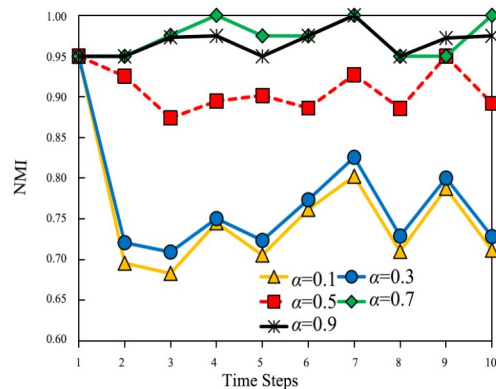
单位:s

| 数据集 | NCI-SeNMF | DNMF | DGR-NMF | FaceNet | ESPRA | Lime |
|---------|-----------|------|---------|-----------|-----------|-----------|
| GN1 | 22 | 20 | 19 | 20 | <u>21</u> | 17 |
| Switch1 | 77 | 78 | 81 | <u>67</u> | 70 | 46 |
| Enron | <u>43</u> | 45 | 52 | 49 | 53 | 25 |
| SFHH | 63 | 65 | 71 | <u>61</u> | 65 | 31 |

注:粗体表示最好结果,下划线表示次优结果.

4.6 α 敏感度分析

NCI-SeNMF 社区发现模型的超参 α 用于平衡网络快照损失项和时变损失项的贡献. 明显地, α 设置为 0 或者 1 时对于动态社区发现都不具有实际意义. 为分析其合理的取值,选择人工合成动态网络 GN1 和真实动态网络 SFHH 作为示例,通过从 0.1~0.9 改变 α 来检测社区,NMI 和 Modularity 分别作为评价指标,最终结果分别如图 6 和图 7 所示. 从图 6 可以看出当 $\alpha \geq 0.7$ 时,所有时刻取得的 NMI 都几乎大于 0.95. 特别地,当 $\alpha = 0.7$ 时,NMI 在 10 个时刻中有 9 个时刻都能够取到最大值,只在时刻 9 时略低于 $\alpha = 0.9$ 取得的 NMI 值. 在图 7 中,也有相似的发现,即当 $\alpha = 0.7$ 时,NCI-SeNMF 在绝大部分时刻可以取得更好的性能. 根据以上观察,在本文的所有实验中, α 都统一设置为 0.7.

图 6 GN1 网络不同 α 值的性能变化

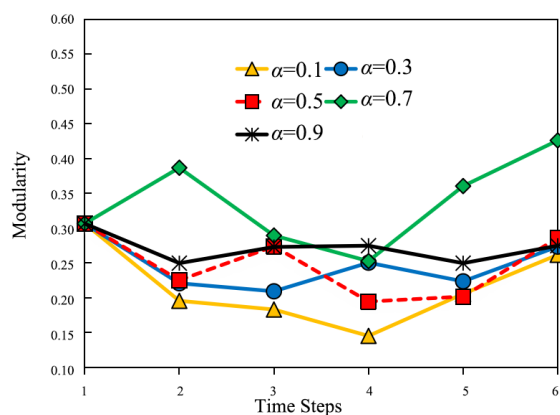


图7 SFHH网络不同 α 值的性能变化

5 结论

动态复杂网络中的前一个时刻网络的社区隶属信息可以用于监督当前时刻网络的社区发现过程,但已有的动态社区发现方法无法有效提取这些信息. 为此本文提出了一种融合节点变化信息的动态社区发现方法 NCI-SeNMF. NCI-SeNMF引入 degeneracy-core 节点和社区转移节点的识别策略,以提取高度可靠的社区隶属关系先验信息,并用于监督当前时刻网络的社区发现过程. 在人工合成动态网络和真实动态网络进行了大量实验分析,结果表明 NCI-SeNMF 具有更稳定且更优的社区发现性能. 虽然 NCI-SeNMF 在小型网络可以获得满意的运行效率,但正如第 3.3 节的分析,NCI-SeNMF 的时间复杂度近似为 $O(n^2)$,这导致它无法有效应用于大规模动态网络. 为此在后续研究中,将重点对节点局部网络结构相似度计算算法和社区发现模型迭代更新算法进行效率优化,目的在于进一步提高 NCI-SeNMF 的整体运行效率.

参考文献

- [1] WANG X F, CHEN G R. Complex networks: Small-world, scale-free and beyond[J]. IEEE Circuits and Systems Magazine, 2003, 3(1): 6-20.
- [2] FORTUNATO S, HRIC D. Community detection in networks: A user guide[J]. Physics Reports, 2016, 659: 1-44.
- [3] MANSOUREH N, MOHAMMAD H F Z, SUSAN B. Fuzzy community detection on the basis of similarities in structural/attribute in large-scale social networks[J]. Artificial Intelligence Review, 2022, 55: 1373-1407.
- [4] LIU S Y, WANG S H. Trajectory community discovery and recommendation by multi-source diffusion modeling[J]. IEEE Transactions on Knowledge and Data Engineering, 2017, 29(4): 898-911.
- [5] HU Z W, LIN A, WILLETT P. Identification of research communities in cited and uncited publications using a co-authorship network[J]. Scientometrics, 2019(118): 1-19.
- [6] DILMAGHANI S, BRUST M R, RIBEIRO C H C, et al. From communities to protein complexes: A local community detection algorithm on PPI networks[J]. PLoS One, 2022, 17(1): e0260484.
- [7] 潘剑飞, 董一鸿, 陈华辉, 等. 基于结构紧密性的重叠社区发现算法[J]. 电子学报, 2019, 47(1): 145-152.
PAN J F, DONG Y H, CHEN H H, et al. The overlapping community discovery algorithm based on compact structure[J]. Acta Electronica Sinica, 2019, 47(1): 145-152. (in Chinese)
- [8] HE D X, LIU H X, FENG Z Y, et al. A joint community detection model: Integrating directed and undirected probabilistic graphical models via factor graph with attention mechanism[J]. IEEE Transactions on Big Data, 2022, 8(4): 994-1006.
- [9] GARZA S E, SCHAEFFER S E. Community detection with the label propagation algorithm: A survey[J]. Physica A: Statistical Mechanics and its Applications, 2019, 534: 122058.
- [10] MOSCATO V, PICARIELLO A, SPERLI G. Community detection based on game theory[J]. Engineering Applications of Artificial Intelligence, 2019, 85: 773-782.
- [11] HE C B, FEI X, CHENG Q W, et al. A survey of community detection in complex networks using nonnegative matrix factorization[J]. IEEE Transactions on Computational Social Systems, 2022, 9(2): 440-457.
- [12] SU X, XUE S, LIU F Z, et al. A comprehensive survey on community detection with deep learning[J]. IEEE Transactions on Neural Networks and Learning Systems, 2024, 35(4): 4682-4702.
- [13] ROSSETTI G, CAZABET R. Community discovery in dynamic networks: A survey[J]. ACM Computing Survey, 2019, 51(2): 35.
- [14] 李赫, 印莹, 李源, 等. 基于多目标演化聚类的大规模动态网络社区检测[J]. 计算机研究与发展, 2019, 56(2): 281-292.
LI H, YIN Y, LI Y, et al. Large-scale dynamic network community detection by multi-objective evolutionary clustering[J]. Journal of Computer Research and Development, 2019, 56(2): 281-292. (in Chinese)
- [15] DHOUIOUI Z, AKAICHI J. Tracking dynamic community evolution in social networks[C]//2014 IEEE/ACM In-

- ternational Conference on Advances in Social Networks Analysis and Mining (ASONAM 2014). Piscataway: IEEE, 2014: 764-770.
- [16] İLHAN N, ÖĞÜDÜCÜ Ş G. Predicting community evolution based on time series modeling[C]//Proceedings of the 2015 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining 2015. New York: ACM, 2015: 1509-1516.
- [17] PENG H, YANG R Y, WANG Z, et al. Lime: Low-cost and incremental learning for dynamic heterogeneous information networks[J]. IEEE Transactions on Computers, 2022, 71(3): 628-642.
- [18] WISSEM I, SABEUR A, HAITHAM M, et al. A distributed and incremental algorithm for large-scale graph clustering[J]. Future Generation Computer Systems, 2022, 34: 334-347.
- [19] AGARWAL P, VERMA R, AGARWAL A, et al. DyPerm: Maximizing permanence for dynamic community detection [C]//PHUNG D, TSENG V, WEBB G, et al. Pacific-Asia Conference on Knowledge Discovery and Data Mining. Cham: Springer, 2018: 437-449.
- [20] YIN Y, ZHAO Y H, LI H, et al. Multi-objective evolutionary clustering for large-scale dynamic community detection[J]. Information Sciences, 2021, 549: 269-287.
- [21] WANG S H, LI G P, HU G Y, et al. Community detection in dynamic networks using constraint non-negative matrix factorization[J]. Intelligent Data Analysis, 2020, 24(1): 119-139.
- [22] MA X K, ZHANG B H, MA C Z, et al. Co-regularized nonnegative matrix factorization for evolving community detection in dynamic networks[J]. Information Sciences, 2020, 528: 265-279.
- [23] SHIN K, ELIASSI-RAD T, FALOUTSOS C. CoreScope: Graph mining using k-core analysis—Patterns, anomalies and algorithms[C]//2016 IEEE 16th International Conference on Data Mining (ICDM). Piscataway: IEEE, 2016: 469-478.
- [24] LANCICHINETTI A, FORTUNATO S, RADICCHI F. Benchmark graphs for testing community detection algorithms[J]. Physical Review E, 2008, 78(4): 046110.
- [25] BLONDEL V D, GUILLAUME J L, LAMBIOTTE R, et al. Fast unfolding of communities in large networks[J]. Journal of Statistical Mechanics: Theory and Experiment, 2008, 10, 10008.
- [26] MALLIAROS F D, VAZIRGIANNIS M. To stay or not to stay: Modeling engagement dynamics in social graphs [C]//Proceedings of the 22nd ACM International Conference on Information & Knowledge Management - CIKM'13. New York: ACM, 2013: 469-478.
- [27] ZHOU T, LÜ L, ZHANG Y C. Predicting missing links via local information[J]. The European Physical Journal B, 2009, 71(4): 623-630.
- [28] 成其伟, 陈启买, 贺超波, 等. 基于改进对称二值非负矩阵分解的重叠社区发现方法[J]. 计算机应用, 2020, 40(11): 3203-3210.
- CHENG Q W, CHEN Q M, HE C B, et al. Overlapping community detection method based on improved symmetric binary nonnegative matrix factorization[J]. Journal of Computer Applications, 2020, 40(11): 3203-3210. (in Chinese)
- [29] JIAO P F, LYU H D, LI X M, et al. Temporal community detection based on symmetric nonnegative matrix factorization[J]. International Journal of Modern Physics B, 2017, 31(13): 1750102.
- [30] LIN Y R, CHI Y, ZHU S H, et al. Facetnet: A framework for analyzing communities and their evolutions in dynamic networks[C]//Proceedings of the 17th International Conference on World Wide Web. New York: ACM, 2008: 685-694.
- [31] WANG P, GAO Z L, MA X K. Dynamic community detection based on network structural perturbation and topological similarity[J]. Journal of Statistical Mechanics: Theory and Experiment, 2017, 2017(1): 013401.
- [32] GREENE D, DOYLE D, CUNNINGHAM P. Tracking the evolution of communities in dynamic social networks [C]//2010 International Conference on Advances in Social Networks Analysis and Mining. Piscataway: IEEE, 2010: 176-183.
- [33] 吴旻诚. 多层网络中结构与动力学研究[D]. 浙江: 浙江大学, 2021.
- WU M C. Research on the Structure and Dynamics in Multilayer Networks[D]. Zhejiang: Zhejiang University,

2021. (in Chinese)

作者简介



贺超波 男, 1981年9月出生于广东省河源市. 现为华南师范大学计算机学院教授、博士生导师, 主要研究方向为图数据挖掘与智能教育.

E-mail: hechaobo@foxmail.com



成其伟 男, 1997年8月出生于广东省江门市. 现工作于维沃移动通信有限公司, 主要研究方向为社区发现.

E-mail: chengqiwei13@163.com



程俊伟 男, 1997年9月出生于福建省宁德市. 现为华南师范大学计算机学院博士研究生, 主要研究方向为图数据挖掘.

E-mail: jung@m.scnu.edu.cn

刘星雨 女, 1998年12月出生于广东省湛江市. 现为华南师范大学计算机学院硕士研究生, 主要研究方向为图数据挖掘.

E-mail: liuxingyu@m.scnu.edu.cn

余 鹏 男, 2000年7月出生于湖北省武汉市. 现为华南师范大学计算机学院硕士研究生, 主要研究方向为图数据挖掘.

E-mail: yupeng@m.scnu.edu.cn

陈启买 男, 1965年5月出生于湖南省衡阳市. 现为华南师范大学教师教学发展中心教授, 主要研究方向为智能教育理论与技术.

E-mail: chenqimai@m.scnu.edu.cn