

动态属性网络的语义社区发现及演化分析方法

贺超波^{1,2}, 成其伟³, 程俊伟¹, 杨佳琦^{1*}, 程 颢¹, 汤 庸^{1,2}

(1. 华南师范大学计算机学院, 广东广州 510631; 2. 琶洲实验室, 广东广州 510335; 3. 维沃移动通信有限公司, 广东东莞 523859)

摘要: 动态属性网络的语义社区发现及演化分析具有重要研究价值, 其包含动态社区发现、社区语义解释及社区演化分析三个任务, 但现有方法均难以同时实现. 针对该问题, 提出一种基于联合非负矩阵分解的方法 DAN-NMF (NMF for Dynamic Attributed Networks). DAN-NMF 可以统一集成网络拓扑结构信息、节点属性信息及社区演化平滑约束信息, 并利用最大最小化优化框架推导相关因子矩阵的迭代更新规则, 从而可以直接获得动态社区发现、社区语义解释及社区演化分析结果. 在人工合成和真实的动态属性网络进行大量相关实验, 结果表明 DAN-NMF 比最优的基准方法在准确性指标上至少提高了 7.3%. 此外, 在真实动态属性网络上的相关数据分析结果也表明 DAN-NMF 能够有效地发现动态社区的演化模式, 并提供丰富的社区语义解释.

关键词: 动态属性网络; 动态社区发现; 社区语义解释; 社区演化分析; 非负矩阵分解

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)11-3757-12

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

DOI: 10.12263/DZXB.20230239

A Community Discovery and Evolution Analysis Method for Dynamic Attributed Networks

HE Chao-bo^{1,2}, CHENG Qi-wei³, CHENG Jun-wei¹, YANG Jia-qi^{1*}, CHENG Hao¹, TANG Yong^{1,2}

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

2. Pazhou Lab, Guangzhou, Guangdong 510335, China;

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

Abstract: The topic of semantic community discovery and evolution analysis in dynamic attributed networks has important research value. It needs to simultaneously accomplish the tasks of dynamic community discovery, community semantic interpretation and community evolution analysis, but existing methods are difficult to achieve this goal. In view of this, this paper proposes a method DAN-NMF (NMF for Dynamic Attributed Networks) based on joint nonnegative matrix factorization. DAN-NMF can uniformly integrate network topology information, attribute information and smooth constraint information from community evolution, and derive iterative update rules of the related factor matrices using the majorization-minimization optimization framework, which helps it to directly obtain the results of dynamic community discovery, community semantic interpretation and community evolution analysis. Extensive experiments are conducted on multiple synthetic and real-world dynamic attributed networks. The results show that DAN-NMF has improved by at least 7.3% in term of accuracy metric, compared to the optimal baseline. Moreover, the data analysis results on real-world dynamic attributed networks also demonstrate that DAN-NMF can effectively discover the evolution patterns of dynamic communities and provide rich community semantic interpretations.

Key words: dynamic attributed networks; dynamic community discovery; community semantic interpretation; community evolution analysis; nonnegative matrix factorization

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

1 引言

现实世界中的各种复杂网络(如社交网络、合著关系网络以及通信网络等)除了具有拓扑结构信息外,还常包含许多额外的属性信息(如节点关联的文本内容特征词、人口统计学属性以及文本标签等).此外,这些网络随着时间的推移,其拓扑结构信息和属性信息都会产生变化,即具备动态性.例如,在社交网络中,用户间的好友关系会动态地建立或消失,即拓扑结构发生改变.此外,如果用户不断发表评论,文本内容的增多会导致其属性信息动态变化.又如在合著关系网络中,作者不断建立合著关系及发表新的论文,也会导致网络拓扑结构及属性信息的变化.形式上,这些具有属性信息,又兼具有动态性的网络,被称为动态属性网络.

动态属性网络包含拓扑结构信息、属性信息以及时间信息,信息更为丰富,可以更好地表示现实世界中的各种复杂交互系统,与静态拓扑网络或者静态属性网络相比也更具有研究价值.语义社区发现及演化分析是动态属性网络的重要研究话题之一,它旨在检测节点间链接紧密、属性高度相似的聚簇,并对聚簇结果进行语义解释和分析其演化关系,即需要同时完成动态社区发现、社区语义解释及社区演化分析三个任务.近几十年来,复杂网络社区发现问题已得到广泛研究,并提出了许多不同类型的解决方法,文献[1~4]分别综述了目前流行的基于标签传播的方法、基于非负矩阵分解(Nonnegative Matrix Factorization, NMF)的方法、基于博弈论的方法以及基于深度学习的方法等.从这些综述文献的相关分析可以看出,尽管已有方法都展示出了一定的有效性,但都最多能同时完成动态社区发现、社区语义解释及社区演化分析中的两个任务,因此动态属性网络的语义社区发现及演化分析仍然需要研究一种更为有效的方法.

基于上述分析,本文提出一种面向动态属性网络的语义社区发现及演化分析方法 DAN-NMF(NMF for Dynamic Attributed Networks). DAN-NMF 首先构建联合 NMF 模型统一集成网络拓扑结构信息、节点属性信息以及社区演化平滑约束信息,然后推导乘性迭代更新规则对模型进行优化求解,最后通过分解得到的相应因子矩阵同时完成动态社区发现、社区语义解释及社区演化分析三个任务.

2 相关工作

2.1 动态社区发现

动态属性网络属于动态网络范畴,目前面向动态网络的社区发现(即动态社区发现)是复杂网络分析的研究热点,现有的相关方法主要包括三类: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.基于两阶段的方法核心思想是第一个阶段先对每一时刻的网络采用静态社区发现算法独立地进行社区划分,然后在第二个阶段对不同时刻的网络社区进行匹配,从而识别社区演化过程,代表性方法包括 CPM (Clique Percolation Method)^[5]、DCDA (Dynamic Community Detection Algorithm)^[6]、ARIMA (AutoRegressive Integrated Moving Average)^[7] 以及 MOCCD (Multi-Objective optimization based on Characteristics for Community Discovery)^[8]等.基于增量学习的方法具有运行效率高的优势,因为其在获得初始社区结构后,只需要利用未来时刻网络的增量变化来改变部分节点的社区隶属情况,而无需再重新考虑整个网络.例如,Agarwal 等人^[9]提出的 DyPerm (Dynamic Permanence)方法、Meng 等人^[10]提出的 IDBLINK (Incremental Density-Based LINK)方法、郭昆等人^[11]提出的 IDCDDC (Incremental Dynamic Community Detection based on Density Clustering)方法以及 Park 等人^[12]提出的 CGC (Contrastive Graph Clustering)方法等都可以高效应用于较大规模的动态网络.基于演化聚类的方法是在假设网络结构变化缓慢条件下,通过演化聚类框架^[13]同时优化两个目标:一是当前时刻网络的社区发现准确性,二是前一时刻与后一时刻网络社区划分结果的相似性.借助统一的优化策略,两个优化目标可以相互增强并提升整体性能.由于演化聚类框架更为可靠且易扩展,大量基于该框架的动态社区发现方法被提出,其中包括 ADFLS (Attention Deep Forest based on Latent Structure)^[14]、DGR-NMF (Dynamic Graph Regularized NMF)^[15]和 Cr-ENMF (Co-regularized Evolutionary NMF)^[16]等方法.

早期的动态社区发现方法更多适用于无属性的动态网络,由于目前动态属性网络日益普遍,一些针对动态属性网络的社区发现方法已开始被陆续提出.例如, Bello 等人^[17]提出对动态属性网络进行社区发现,采用融合节点属性信息的 Louvain 方法^[18]对每个时刻的网络进行社区划分. Appel 等人^[19]提出同时考虑链接、内容及时间信息的共享因子分解模型 Chimera,以多维形式提取网络的潜在语义社区结构. Zhang 等人^[20]提出集成网络拓扑结构、文本及时间信息的生成模型 DTCD (Dynamic Topical Community Detection),可以发现结构和主题上都有意义的社区.此外, Li 等人^[21]通过将时间与癌症基因属性信息集成到 TANMF (Temporal Attributed NMF)模型,能够准确检测癌症动态属性网络存在的社区.

总的来说,现有动态社区发现方法最多能同时完成动态属性网络语义社区发现及演化分析涉及的其中两个任务.例如, DGR-NMF 及 Cr-ENMF 只关注于利用拓扑结构信息进行动态社区发现. DTCD、Chimera 及

TANMF等虽然集成了拓扑信息和属性信息进行动态社区发现及社区语义解释,但并没有考虑社区演化分析,CPM、DCDA及ARIMA等虽然能检测动态社区和追踪社区演化,但没有集成利用属性信息进行社区语义解释.区别于已有的方法,本文提出的DAN-NMF方法能够同时完成动态社区发现、社区语义解释及社区演化分析三个任务.

2.2 NMF与社区发现

NMF是一种包含非负值约束的低秩矩阵分解方法^[22].形式上,给定一个大小为 $d_1 \times d_2$ 的非负矩阵 Y 和一个正整数 $k(k \ll \min(d_1, d_2))$,NMF可以将 Y 分解为两个大小分别为 $d_1 \times k, k \times d_2$ 的非负矩阵 P 和 Q ,且 P 和 Q 的乘积接近于 $Y: Y \approx PQ$,常用Frobenius范数的平方 $\|Y - PQ\|_F^2$ 度量其中误差.文献[23]指出,当 $d_1 = d_2$,即是对称矩阵时,NMF可以进一步转换为与聚类算法K-means等价的对称分解形式: $Y \approx PP^T$,其中 P 可以视为聚类结果指示矩阵.由于对称NMF具有聚类能力,目前已被大量用于解决复杂网络的社区发现问题,此时 Y 和 P 分别作为拓扑结构矩阵和社区隶属强度表示矩阵, PP^T 被视为基于两两节点的社区隶属强度表示向量重建网络的拓扑结构 Y .文献[2]较为全面地归纳分析了基于对称NMF的各类社区发现方法,并指出对称NMF具有很好的可扩展性,可以通过联合对称NMF或者加入正则约束项处理各种网络的社区发现问题.例如,DGR-NMF^[15]、Cr-ENMF^[16]及RSECD^[24]等均是对称NMF的扩展方法.本文提出的动态属性网络社区发现方法DAN-NMF也是基于对称NMF进行扩展,通过加入联合分解项及正则项集成更多相关信息以获得更为准确的动态社区发现结果.

3 语义社区发现及演化分析方法 DAN-NMF

3.1 问题定义与符号说明

对于给定的动态属性网络,本文将其建模为一个包含 T 个时刻的快照属性网络集合,表示为 $G = \{G_1, G_2, \dots, G_T\}$,其中任何一个快照属性网络都建模为无向无权图 $G = (V, E, X)$,其中 $V = \{v_1, v_2, \dots, v_n\}$ 为节点

集合, $E = \{e_1, e_2, \dots, e_l\}$ 为边集合, $X = \{x_1, x_2, \dots, x_m\}$ 为节点属性集合.对于节点集合 V 中的任意一个节点 v_i ,都关联一个 $1 \times m$ 维的属性向量 $S_i = [s_{i1}, s_{i2}, \dots, s_{im}]$,其元素 s_{ij} 取值为0或1,取值为1时表示节点 v_i 具备属性 x_j ,取值为0时则相反.对于给定的一个动态属性网络 G ,本文的问题在于同时完成其语义社区发现及演化分析涉及三个任务:动态社区发现、社区语义解释以及社区演化分析,这三个任务的具体要求分别说明如下:

(1)动态社区发现:给定一个动态属性网络 $G = \{G_1, G_2, \dots, G_T\}$,将其每个时刻的属性网络 $G_t \in G$ 划分为 k 个不相交的社区 $C_t = \{C_{1,t}, C_{2,t}, \dots, C_{k,t}\}$,其中对于 G_t 中的任意两个社区 $C_{i,t}, C_{j,t} \in C_t$ 都有 $C_{i,t} \cap C_{j,t} = \emptyset$.此外,同一社区节点之间的链接相对于不同社区节点之间的链接应更为稠密,同一社区节点间所具有的节点属性相对于不同社区节点之间的节点属性也应更为相似.

(2)社区语义解释:给定第 t 时刻的第 i 个社区 $C_{i,t}$,需要推断该社区 $C_{i,t}$ 所对应的潜在主题 $Z_{i,t}$.主题 $Z_{i,t}$ 可以使用一个向量 $(Z_{1,i,t}, Z_{2,i,t}, \dots, Z_{m,i,t})$ 表示,其中 $Z_{j,i,t}$ 表示属性 j 隶属于主题 $Z_{i,t}$ 的强度.

(3)社区演化关系分析:给定相邻时刻的任意两个社区 $C_{i,t-1}$ 和 $C_{j,t}$,可以量化社区 $C_{i,t-1}$ 转移到社区 $C_{j,t}$ 的概率 p_{ij} .

为方便后续介绍DAN-NMF方法的具体实现,在表1列出常用符号.

3.2 模型构建

(1)建模拓扑结构信息.对任一给定的第 t 时刻快照属性网络 G_t ,定义一个表示节点社区隶属强度的矩阵 H_t ,其元素 $H_{ir,t}$ 表示为在时刻 t 中,节点 i 隶属于社区 r 的强度. $H_{ir,t}$ 与 $H_{jr,t}$ 的乘积可以视为节点 i 与节点 j 在第 r 个社区中的交互概率.通过对所有的 k_t 个社区进行求和: $\sum_{r=1}^{k_t} H_{ir,t} H_{jr,t}$,可得到节点 i 与节点 j 总的交互概率.这种节点预期交互概率的计算过程表明如果两个节点具有相似的社区隶属关系则它们之间将有可能建立

表1 符号说明

| 符号 | 说明 |
|-------|--|
| A_t | 第 t 时刻属性网络快照的邻接矩阵,用于表示属性网络快照的拓扑结构信息,大小为 $n \times n$ |
| S_t | 第 t 时刻属性网络快照的节点属性矩阵,大小为 $n \times m$ |
| k_t | 第 t 时刻属性网络快照的社区数和主题数 |
| H_t | 第 t 时刻节点社区隶属强度矩阵,大小为 $n \times k_t$ |
| B_t | 第 t 时刻节点主题隶属强度矩阵,大小为 $n \times k_t$ |
| W_t | 第 t 时刻属性主题隶属强度矩阵,大小为 $m \times k_t$ |
| M_t | 第 t 时刻社区主题映射矩阵,大小为 $k_t \times k_t$ |
| G_t | 第 $t-1$ 时刻到第 t 时刻的社区演化关系矩阵,大小为 $k_{t-1} \times k_t$ |

边, 即有 $A_{ij,t} \approx \sum_{r=1}^{k_i} H_{ir,t} H_{jr,t}$. 将该式扩展到所有节点对, 可建模通过社区重建网络拓扑结构的损失为

$$L_{\text{link}}(\mathbf{H}_t) = \sum_{i=1}^n \sum_{j=1}^n \left(A_{ij,t} - \sum_{r=1}^{k_i} H_{ir,t} H_{jr,t} \right)^2 \quad (1)$$

$$= \|\mathbf{A}_t - \mathbf{H}_t \mathbf{H}_t^T\|_F^2$$

(2) 建模属性信息. 令 \mathbf{B}_t 表示节点主题隶属强度矩阵, 其元素 $B_{ip,t}$ 的含义为在时刻 t 中, 节点 i 隶属于主题 p 的强度. 令 \mathbf{W}_t 表示属性主题隶属强度矩阵, 其元素 $W_{jp,t}$ 的含义为在时刻 t 中, 属性 j 隶属于主题 p 的强度. 节点 i 具备属性 j 的概率可通过 $\sum_{p=1}^{k_j} B_{ip,t} W_{jp,t}$ 进行计算, 该计算过程表明若节点 i 与属性 j 隶属于同一个潜在主题, 那么节点 i 应倾向具备属性 j , 即有 $S_{ij,t} \approx \sum_{p=1}^{k_j} B_{ip,t} W_{jp,t}$, 扩展至所有的节点属性对, 可建立如下损失项:

$$L_{\text{attribute}}(\mathbf{B}_t, \mathbf{W}_t) = \sum_{i=1}^n \sum_{j=1}^m \left(S_{ij,t} - \sum_{p=1}^{k_j} B_{ip,t} W_{jp,t} \right)^2 \quad (2)$$

$$= \|\mathbf{S}_t - \mathbf{B}_t \mathbf{W}_t^T\|_F^2$$

(3) 建模拓扑结构信息和属性信息的融合. 为了促进拓扑信息和属性信息相互融合, 同时令模型能够具备一定的鲁棒性, 以便应对社区与主题相互冲突的情况, 引入社区主题映射对称矩阵 \mathbf{M}_t , 其元素 $M_{ij,t}$ 表明在时刻 t 中, 社区 i 映射到主题 j 的强度.

具体损失项的构造如下:

$$L_{\text{combine}}(\mathbf{H}_t, \mathbf{B}_t) = \|\mathbf{H}_t \mathbf{M}_t - \mathbf{B}_t\|_F^2 + \|\mathbf{M}_t - \mathbf{I}\|_F^2 \quad (3)$$

其中, \mathbf{I} 为单位矩阵, 用于约束每一个社区都尽可能地对应一个主题.

(4) 建模社区演化模式. 引入社区转移矩阵 \mathbf{G}_t 描述第 $t-1$ 时刻到第 t 时刻的社区演化模式, 其元素 $G_{ij,t}$ 表示第 $t-1$ 时刻的社区 i 在第 t 时刻转移到社区 j 的概率. 当 \mathbf{H}_{t-1} 与 \mathbf{G}_t 的乘积近似于 \mathbf{H}_t 时, 即 $\mathbf{H}_{t-1} \mathbf{G}_t \approx \mathbf{H}_t$, \mathbf{G}_t 能表达相邻时刻社区间的转移概率, 由此产生的损失可表示为

$$L_{\text{evolution}}(\mathbf{H}_t) = \|\mathbf{H}_{t-1} \mathbf{G}_t - \mathbf{H}_t\|_F^2 \quad (4)$$

(5) 建模社区演化的平滑性. 动态属性网络的节点社区隶属关系通常具有平滑性, 即节点在相邻时刻的社区隶属关系变化是缓慢的. 为保留该平滑性, 一般地在 $t-1$ 时刻的属性网络中两个节点的局部拓扑结构相似性越大, 则认为在当前时刻 t 中这两个节点越可能隶属于同一个社区. 为提高计算效率, 这里直接使用上一时刻表示节点间是否连接的 0/1 值作为相似性, 由此可构造如下损失项:

$$L_{\text{smooth}_C}(\mathbf{H}_t) = \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n A_{ij,t-1} \|\mathbf{H}_{i,t-1} - \mathbf{H}_{j,t-1}\|_F^2 \quad (5)$$

$$= \text{tr}(\mathbf{H}_t^T (D(\mathbf{A}_{t-1}) - \mathbf{A}_{t-1}) \mathbf{H}_t)$$

其中, $D(\mathbf{X})$ 以矩阵 \mathbf{X} 为输入, 输出如下对角矩阵:

$$D(\mathbf{X})_{ij} = \begin{cases} \sum_{p=1}^n X_{ip}, & \text{if } i=j \\ 0, & \text{if } i \neq j \end{cases} \quad (6)$$

同样地, 节点主题隶属关系也具有平滑性, 即在相邻时刻节点所属主题的变化是缓慢的. 可以认为在 $t-1$ 时刻的属性网络中, 如果两个节点在主题状态空间中距离越近, 那么在当前时刻中这两个节点越可能隶属于同一社区, 对应损失项为

$$L_{\text{smooth}_T}(\mathbf{H}_t) = \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n (B_{i,t-1} B_{j,t-1}^T)_{ij} \|\mathbf{H}_{i,t-1} - \mathbf{H}_{j,t-1}\|_F^2 \quad (7)$$

$$= \text{tr}(\mathbf{H}_t^T (D(\mathbf{B}_{t-1} \mathbf{B}_{t-1}^T) - \mathbf{B}_{t-1} \mathbf{B}_{t-1}^T) \mathbf{H}_t)$$

(6) 统一模型. 以上损失项涉及到的各个矩阵都具有非负值约束, 因此可以基于联合 NMF 框架集成上述损失项构建如下面向动态属性网络的语义社区发现及演化分析模型:

$\min L(\mathbf{H}_t, \mathbf{B}_t, \mathbf{W}_t, \mathbf{M}_t, \mathbf{G}_t)$

$$= \begin{cases} \|\mathbf{A}_t - \mathbf{H}_t \mathbf{H}_t^T\|_F^2 + \|\mathbf{S}_t - \mathbf{B}_t \mathbf{W}_t^T\|_F^2 \\ \quad + \alpha (\|\mathbf{H}_t \mathbf{M}_t - \mathbf{B}_t\|_F^2 + \|\mathbf{M}_t - \mathbf{I}\|_F^2), & \text{if } t=1 \\ \|\mathbf{A}_t - \mathbf{H}_t \mathbf{H}_t^T\|_F^2 + \|\mathbf{S}_t - \mathbf{B}_t \mathbf{W}_t^T\|_F^2 \\ \quad + \alpha (\|\mathbf{H}_t \mathbf{M}_t - \mathbf{B}_t\|_F^2 + \|\mathbf{M}_t - \mathbf{I}\|_F^2) \\ \quad + \beta \|\mathbf{H}_{t-1} \mathbf{G}_t - \mathbf{H}_t\|_F^2 + \gamma \text{tr}(\mathbf{H}_t^T (F(\mathbf{A}_{t-1})) \mathbf{H}_t) & \text{if } t > 1 \\ \quad + \eta \text{tr}(\mathbf{H}_t^T (F(\mathbf{B}_{t-1} \mathbf{B}_{t-1}^T)) \mathbf{H}_t), & \end{cases} \quad (8)$$

其中, $\alpha, \beta, \gamma, \eta$ 为相应损失项的权重参数, $F(\mathbf{X}) = D(\mathbf{X}) - \mathbf{X}$. 当 $t=1$ 时, 动态属性网络只有 1 个快照, 因此不需要建模社区演化模式及平滑性.

3.3 模型优化及算法分析

采用最大最小化框架^[25]可以迭代优化求解式(8)目标函数的局部最优解. 该框架轮流迭代更新变量 \mathbf{H}_t 、 \mathbf{B}_t 、 \mathbf{W}_t 、 \mathbf{M}_t 及 \mathbf{G}_t , 在更新其中一个变量时, 将其余变量看作常数, 因此可以将式(8)模型的优化问题划分为 5 个变量的子优化问题. 具体地, 由于每个变量都具有非负值约束, 那么对每个变量的优化都可以视为受限求极值问题, 可以通过拉格朗日乘数法和 KKT (Karush-Kuhn-Tucker) 条件推导其迭代更新规则, 下面分别给出每个变量的优化过程.

(1) 优化变量 \mathbf{H}_t

当 $t > 1$ 时, 将 \mathbf{B}_t 、 \mathbf{W}_t 、 \mathbf{M}_t 、 \mathbf{G}_t 看作常量. 由于 $\|\mathbf{X}\|_F^2 =$

$\text{tr}(\mathbf{X}\mathbf{X}^T), \text{tr}(\mathbf{X}) = \text{tr}(\mathbf{X}^T)$, 式(8)可被重写为

$$\begin{aligned} L(\mathbf{H}_t) = & \text{tr}(\mathbf{H}_t \mathbf{H}_t^T \mathbf{H}_t \mathbf{H}_t^T - 2\mathbf{A}_t \mathbf{H}_t \mathbf{H}_t^T) \\ & + \alpha \text{tr}(\mathbf{M}_t^T \mathbf{H}_t^T \mathbf{H}_t \mathbf{M}_t - 2\mathbf{M}_t^T \mathbf{H}_t^T \mathbf{B}_t) \\ & + \beta \text{tr}(\mathbf{H}_t \mathbf{H}_t^T - 2\mathbf{H}_t \mathbf{G}_t^T \mathbf{H}_t^T) \\ & + \gamma \text{tr}(\mathbf{H}_t^T (D(\mathbf{A}_{t-1}) - \mathbf{A}_{t-1}) \mathbf{H}_t) \\ & + \eta \text{tr}(\mathbf{H}_t^T (D(\mathbf{B}_{t-1} \mathbf{B}_{t-1}^T) - \mathbf{B}_{t-1} \mathbf{B}_{t-1}^T) \mathbf{H}_t) \end{aligned} \quad (9)$$

设 Ψ 为约束 $\mathbf{H}_t \geq 0$ 的拉格朗日乘子, 式(9)的拉格朗日函数为

$$L(\mathbf{H}_t) = L(\mathbf{H}_t) + \text{tr}(\Psi \mathbf{H}_t^T) \quad (10)$$

$L(\mathbf{H}_t)$ 对 \mathbf{H}_t 求导可得:

$$H_{ij,t} = H_{ij,t} \frac{(2\mathbf{A}_t \mathbf{H}_t + \alpha \mathbf{B}_t \mathbf{M}_t^T + \beta \mathbf{H}_{t-1} \mathbf{G}_t + \gamma \mathbf{A}_{t-1} \mathbf{H}_t + \eta \mathbf{B}_{t-1} \mathbf{B}_{t-1}^T \mathbf{H}_t)_{ij}}{(2\mathbf{H}_t \mathbf{H}_t^T \mathbf{H}_t + \alpha \mathbf{H}_t \mathbf{M}_t \mathbf{M}_t^T + \beta \mathbf{H}_t + \gamma D(\mathbf{A}_{t-1}) \mathbf{H}_t + \eta D(\mathbf{B}_{t-1} \mathbf{B}_{t-1}^T) \mathbf{H}_t)_{ij}} \quad (13)$$

当 $t=1$ 时, 由上述过程同理可得 $H_{ij,t}$ 的乘性迭代更新规则:

$$H_{ij,t} = H_{ij,t} \frac{(2\mathbf{A}_t \mathbf{H}_t + \alpha \mathbf{B}_t \mathbf{M}_t^T)_{ij}}{(2\mathbf{H}_t \mathbf{H}_t^T \mathbf{H}_t + \alpha \mathbf{H}_t \mathbf{M}_t \mathbf{M}_t^T)_{ij}} \quad (14)$$

(2) 优化变量 \mathbf{B}_t

将 $\mathbf{H}_t, \mathbf{W}_t, \mathbf{M}_t, \mathbf{G}_t$ 看作常量, 式(8)在 $t=1$ 和 $t>1$ 时是等价的, 其可被重写为

$$\begin{aligned} L(\mathbf{B}_t) = & \text{tr}(\mathbf{W}_t \mathbf{B}_t^T \mathbf{B}_t \mathbf{W}_t^T - 2\mathbf{S}_t^T \mathbf{B}_t \mathbf{W}_t^T) \\ & + \alpha \text{tr}(\mathbf{B}_t^T \mathbf{B}_t - 2\mathbf{M}_t^T \mathbf{H}_t^T \mathbf{B}_t) \end{aligned} \quad (15)$$

设 Ψ 为约束 $\mathbf{B}_t \geq 0$ 的拉格朗日乘子, 式(15)的拉格朗日函数为

$$L(\mathbf{B}_t) = L(\mathbf{B}_t) + \text{tr}(\Psi \mathbf{B}_t^T) \quad (16)$$

$L(\mathbf{B}_t)$ 对 \mathbf{B}_t 求导可得

$$\frac{\partial L(\mathbf{B}_t)}{\partial \mathbf{B}_t} = \mathbf{B}_t \mathbf{W}_t^T \mathbf{W}_t - \mathbf{S}_t \mathbf{W}_t + \alpha(\mathbf{B}_t - \mathbf{H}_t \mathbf{M}_t) + \Psi \quad (17)$$

根据 KKT 条件 $\psi_{ij} B_{ij,t} = 0$ 可得

$$\left[(\mathbf{B}_t \mathbf{W}_t^T \mathbf{W}_t)_{ij} - (\mathbf{S}_t \mathbf{W}_t)_{ij} + \alpha(\mathbf{B}_t - \mathbf{H}_t \mathbf{M}_t)_{ij} \right] B_{ij,t} = 0 \quad (18)$$

由式(18)可得 $B_{ij,t}$ 的乘性迭代更新规则为

$$B_{ij,t} = B_{ij,t} \frac{(\mathbf{S}_t \mathbf{W}_t + \alpha \mathbf{H}_t \mathbf{M}_t)_{ij}}{(\mathbf{B}_t \mathbf{W}_t^T \mathbf{W}_t + \alpha \mathbf{B}_t)_{ij}} \quad (19)$$

(3) 优化变量 \mathbf{W}_t

将 $\mathbf{H}_t, \mathbf{B}_t, \mathbf{M}_t, \mathbf{G}_t$ 看作常量, 式(8)在 $t=1$ 和 $t>1$ 时是等价的, 其可被重写为

$$L(\mathbf{W}_t) = \text{tr}(\mathbf{W}_t \mathbf{B}_t^T \mathbf{B}_t \mathbf{W}_t^T - 2\mathbf{S}_t^T \mathbf{B}_t \mathbf{W}_t^T) \quad (20)$$

设 Ψ 为约束 $\mathbf{W}_t \geq 0$ 的拉格朗日乘子, 式(20)的拉格朗日函数为

$$L(\mathbf{W}_t) = L(\mathbf{W}_t) + \text{tr}(\Psi \mathbf{W}_t^T) \quad (21)$$

$$\begin{aligned} \frac{\partial L(\mathbf{H}_t)}{\partial \mathbf{H}_t} = & 2\mathbf{H}_t \mathbf{H}_t^T \mathbf{H}_t - 2\mathbf{A}_t \mathbf{H}_t + \alpha(\mathbf{H}_t \mathbf{M}_t \mathbf{M}_t^T - \mathbf{B}_t \mathbf{M}_t^T) \\ & + \beta(\mathbf{H}_t - \mathbf{H}_{t-1} \mathbf{G}_t) + \gamma(D(\mathbf{A}_{t-1}) - \mathbf{A}_{t-1}) \mathbf{H}_t \\ & + \eta(D(\mathbf{B}_{t-1} \mathbf{B}_{t-1}^T) - \mathbf{B}_{t-1} \mathbf{B}_{t-1}^T) \mathbf{H}_t + \Psi \end{aligned} \quad (11)$$

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

$$\begin{aligned} & \left[2(\mathbf{H}_t \mathbf{H}_t^T \mathbf{H}_t)_{ij} - 2(\mathbf{A}_t \mathbf{H}_t)_{ij} + \alpha(\mathbf{H}_t \mathbf{M}_t \mathbf{M}_t^T - \mathbf{B}_t \mathbf{M}_t^T)_{ij} \right. \\ & \left. + \beta(\mathbf{H}_t - \mathbf{H}_{t-1} \mathbf{G}_t)_{ij} + \gamma((D(\mathbf{A}_{t-1}) - \mathbf{A}_{t-1}) \mathbf{H}_t)_{ij} \right. \\ & \left. + \eta((D(\mathbf{B}_{t-1} \mathbf{B}_{t-1}^T) - \mathbf{B}_{t-1} \mathbf{B}_{t-1}^T) \mathbf{H}_t)_{ij} \right] H_{ij,t} = 0 \end{aligned} \quad (12)$$

由式(12)可得 $H_{ij,t}$ 的乘性迭代更新规则:

$L(\mathbf{W}_t)$ 对 \mathbf{W}_t 求导可得

$$\frac{\partial L(\mathbf{W}_t)}{\partial \mathbf{W}_t} = \mathbf{W}_t \mathbf{B}_t^T \mathbf{B}_t - \mathbf{S}_t \mathbf{B}_t + \Psi \quad (22)$$

根据 KKT 条件 $\psi_{ij} W_{ij,t} = 0$, 可得

$$\left[(\mathbf{W}_t \mathbf{B}_t^T \mathbf{B}_t)_{ij} - (\mathbf{S}_t \mathbf{B}_t)_{ij} \right] W_{ij,t} = 0 \quad (23)$$

由式(23)可得 $W_{ij,t}$ 的乘性迭代更新规则为

$$W_{ij,t} = W_{ij,t} \frac{(\mathbf{S}_t \mathbf{B}_t)_{ij}}{(\mathbf{W}_t \mathbf{B}_t^T \mathbf{B}_t)_{ij}} \quad (24)$$

(4) 优化变量 \mathbf{M}_t

将 $\mathbf{H}_t, \mathbf{B}_t, \mathbf{W}_t, \mathbf{G}_t$ 看作常量, 式(8)在 $t=1$ 和 $t>1$ 时是等价的, 其可被重写为

$$\begin{aligned} L(\mathbf{M}_t) = & \text{tr}(\mathbf{M}_t^T \mathbf{H}_t^T \mathbf{H}_t \mathbf{M}_t - 2\mathbf{M}_t^T \mathbf{H}_t^T \mathbf{B}_t) \\ & + \text{tr}(\mathbf{M}_t^T \mathbf{M}_t - 2\mathbf{M}_t^T) \end{aligned} \quad (25)$$

设 Ψ 为约束 $\mathbf{M}_t \geq 0$ 的拉格朗日乘子, 式(25)的拉格朗日函数如下:

$$L(\mathbf{M}_t) = L(\mathbf{M}_t) + \text{tr}(\Psi \mathbf{M}_t^T) \quad (26)$$

$L(\mathbf{M}_t)$ 对 \mathbf{M}_t 求导可得到式(27):

$$\frac{\partial L(\mathbf{M}_t)}{\partial \mathbf{M}_t} = \mathbf{H}_t^T \mathbf{H}_t \mathbf{M}_t - \mathbf{H}_t^T \mathbf{B}_t + \mathbf{M}_t - \mathbf{I} + \Psi \quad (27)$$

根据 KKT 条件 $\psi_{ij} M_{ij,t} = 0$, 可得

$$\left[(\mathbf{H}_t^T \mathbf{H}_t \mathbf{M}_t)_{ij} - (\mathbf{H}_t^T \mathbf{B}_t)_{ij} + M_{ij,t} + I_{ij} \right] M_{ij,t} = 0 \quad (28)$$

由上式可得 $W_{ij,t}$ 的乘性迭代更新规则为

$$M_{ij,t} = M_{ij,t} \frac{(\mathbf{H}_t^T \mathbf{B}_t + \mathbf{I})_{ij}}{(\mathbf{H}_t^T \mathbf{H}_t \mathbf{M}_t + \mathbf{M}_t)_{ij}} \quad (29)$$

(5) 优化变量 \mathbf{G}_t

由于 \mathbf{G}_t 在 $t=1$ 时不存在, 因此只需考虑 $t>1$ 的情

况. 将 H_t, B_t, W_t, M_t 看作常量, 式(8)可被重写为

$$L(\mathbf{G}_t) = \text{tr}(\mathbf{G}_t^T \mathbf{H}_{t-1}^T \mathbf{H}_{t-1} \mathbf{G}_t - 2\mathbf{G}_t^T \mathbf{H}_{t-1}^T \mathbf{H}_t) \quad (30)$$

设 Ψ 为约束 $\mathbf{G}_t \geq 0$ 的拉格朗日乘子, 式(30)的拉格朗日函数如下:

$$L(\mathbf{G}_t) = L(\mathbf{G}_t) + \text{tr}(\Psi \mathbf{G}_t^T) \quad (31)$$

$L(\mathbf{G}_t)$ 对 \mathbf{G}_t 求导可得到式(32):

$$\frac{\partial L(\mathbf{G}_t)}{\partial \mathbf{G}_t} = \mathbf{H}_{t-1}^T \mathbf{H}_{t-1} \mathbf{G}_t - \mathbf{H}_{t-1}^T \mathbf{H}_t + \Psi \quad (32)$$

根据 KKT 条件 $\psi_{ij} G_{ij,t} = 0$, 可得

$$\left[(\mathbf{H}_{t-1}^T \mathbf{H}_{t-1} \mathbf{G}_t)_{ij} - (\mathbf{H}_{t-1}^T \mathbf{H}_t)_{ij} \right] G_{ij,t} = 0 \quad (33)$$

由(33)可得 $G_{ij,t}$ 的乘性迭代更新规则为

$$G_{ij,t} = G_{ij,t} \frac{(\mathbf{H}_{t-1}^T \mathbf{H}_t)_{ij}}{(\mathbf{H}_{t-1}^T \mathbf{H}_{t-1} \mathbf{G}_t)_{ij}} \quad (34)$$

使用 DAN-NMF 方法可对动态属性网络同时进行动态社区发现、社区语义解释及社区演化关系分析, 其总的执行流程如算法 1 所示.

从算法 1 可以看出, DAN-NMF 最耗时的部分来自于各个矩阵变量的迭代更新操作. 具体地, 迭代更新 H_t, M_t 及 G_t 的时间复杂度为 $O(kn^2)$, 而迭代更新 B_t 和 W_t 的时间复杂度为 $O(knm)$, 其中 k 为所有属性网络快照的平均社区个数. 假设每个属性网络快照的平均迭代更新次数为 t_{iter} , 则算法总的的时间复杂度为 $O(t_{\text{iter}} Tkn(3n+2m))$.

4 实验

4.1 数据集

(1) 人工合成动态属性网络

目前人工合成动态网络有两种较为流行的生成工具: Dynamic LFR^[26]和 RDyn^[27], 但两种工具生成的网络都没有属性信息, 为此需要为 Dynamic LFR 和 RDyn 的合成动态网络统一定义节点属性信息的生成过程. 首先为每个节点生成属于该节点的 $1 \times m$ 维属性向量, 其中属性向量生成的规则为: 假若节点的真实社区编号为 x (编号从 1 开始), 属性向量下标从 $(x-1) \times h$ 至 $x \times h - 1$ 按照参数 $p_{\text{in}} \in [0, 1]$ 的概率生成 1, 剩余的位置按照参数 $p_{\text{out}} \in [0, 1]$ 的概率生成 0. 最后随机交换网络 $p_{\text{mis}} \times 100\%$ 节点的属性向量, 其中 $p_{\text{mis}} \in [0, 1]$, 这可以模拟属性信息存在噪声的情况. 接下来分别介绍 Dynamic LFR 和 RDyn 生成的动态属性网络.

(a) Dynamic LFR 网络. Dynamic LFR 根据网络演化事件的不同生成 4 组共 8 个网络, 每个网络的参数设置如表 2 所示, 其中 n 表示网络节点总个数, T 表示网络快照个数, l_{max} 为各快照网络的最大边数, μ 表示混合参数. 参

算法 1 DAN-NMF 语义社区发现及演化分析

输入: T 个时刻的属性网络邻接矩阵 $\mathbf{A}_s = (\mathbf{A}_1, \mathbf{A}_2, \dots, \mathbf{A}_T)$ 、节点属性矩阵 $\mathbf{S}_s = (\mathbf{S}_1, \mathbf{S}_2, \dots, \mathbf{S}_T)$ 、社区数 $k_s = \{k_1, k_2, \dots, k_T\}$ 、权重参数 $\alpha, \beta, \gamma, \eta$

输出: T 个时刻的社区划分结果 $C_s = (C_1, C_2, \dots, C_T)$ 、社区对应主题

$\mathbf{Z}_s = (\mathbf{Z}_1, \mathbf{Z}_2, \dots, \mathbf{Z}_T)$ 、 $T-1$ 个社区演化矩阵 $\mathbf{G}_s = (\mathbf{G}_2, \mathbf{G}_3, \dots, \mathbf{G}_T)$

1: $C_s = ()$ // 保存社区划分结果的矩阵

2: $\mathbf{Z}_s = ()$ // 保存社区对应主题的矩阵

3: $\mathbf{G}_s = ()$ // 保存社区演化矩阵

4: FOR $t = 1$ TO T

5: 随机初始化 H_t, B_t, W_t, M_t, G_t

6: WHILE 目标函数式(8)不收敛 DO

7: IF $t = 1$

8: 根据式(14)更新 H_t

9: ELSE

10: 根据式(13)更新 H_t

11: 根据式(34)更新 G_t

12: END IF

13: 根据式(19)更新 B_t

14: 根据式(24)更新 W_t

15: 根据式(29)更新 M_t

16: END WHILE

17: 根据 H_t 获得时刻 t 的社区划分结果 C_t , 添加到 C_s 中

18: IF $t > 1$

19: 对 G_t 进行归一化, 即 $\forall i, \sum_{t=1}^n G_{i,t} = 1$, 然后添加到 \mathbf{G}_s

20: END IF

21: // 获取每个社区对应的主题

22: FOR $C_{i,t}$ IN C_t

23: $p = \arg \max_i M_{ii}$

24: $\mathbf{Z}_{i,t} = \mathbf{W}_p$

25: END FOR

26: 将 \mathbf{Z}_t 添加到集合 \mathbf{Z}_s 中

27: END FOR

28: RETURN $C_s, \mathbf{Z}_s, \mathbf{G}_s$

数 p 在各个网络的含义有所不同, 在 Switch 网络中 p 表示的是节点相邻时刻发生社区转移的概率. 在 BirthDeath 网络中 p 表示的是: 每个时刻中, $p \times 100\%$ 的社区发生出生事件及 $p \times 100\%$ 的社区发生死亡事件; 在 ExpandContract 网络中 p 表示的是: 每个时刻中, $p \times 100\%$ 的社区会发生扩张事件及 $p \times 100\%$ 的社区发生收缩事件; 在 MergeSplit 网络中 p 表示的是: 每个时刻中, $p \times 100\%$ 的社区发生合并事件及 $p \times 100\%$ 的社区发生分裂事件. 参数 r 则是 ExpandContract 网络特有的参数, 表示社区扩张或收缩的比例. 通过表 2 的参数设置情况可以看出, 对于同组的网络, 编号为 2 的网络比编号为 1 的网络进行社区发现的难度要更大, 因为编号为 2 的网络相比编号为 1 的网络: 演

表 2 Dynamic LFR 网络参数设置

| 网络名称 | n | T | l_{\max} | μ | P | r | P_{in} | P_{out} | P_{mis} | h |
|----------------|------|-----|------------|-------|-------|-----|----------|-----------|-----------|-----|
| Switch1 | 400 | 10 | 1524 | 0.3 | 0.2 | — | 0.5 | 0.05 | 0.1 | 20 |
| Switch2 | 400 | 10 | 1927 | 0.3 | 0.4 | — | 0.5 | 0.1 | 0.2 | 20 |
| BirthDeath1 | 400 | 10 | 1350 | 0.3 | 0.125 | — | 0.5 | 0.05 | 0.1 | 20 |
| BirthDeath2 | 400 | 10 | 1505 | 0.3 | 0.25 | — | 0.5 | 0.1 | 0.2 | 20 |
| MergeSplit1 | 400 | 10 | 1491 | 0.3 | 0.25 | — | 0.5 | 0.05 | 0.1 | 20 |
| MergeSplit2 | 400 | 10 | 1526 | 0.3 | 0.25 | — | 0.5 | 0.1 | 0.2 | 20 |
| ExpandContrat1 | 5000 | 10 | 51718 | 0.3 | 0.125 | 0.2 | 0.5 | 0.05 | 0.1 | 20 |
| ExpandContrat2 | 5000 | 10 | 50404 | 0.3 | 0.25 | 0.4 | 0.5 | 0.1 | 0.2 | 20 |

化事件的发生更为频繁,节点属性信息存在更多的噪声以及属性信息与社区结构的匹配程度更低.

(b)RDyn 网络. 动态网络生成器 RDyn 不仅能模拟交互动态,还能模拟社区动态,这里利用 RDyn 一共生成 4 个网络: Network1、Network2、Network3 以及 Network4,各网络的参数如表 3 所示,其中 iter 为 RDyn 生成网络过程的迭代次数, q 表示电导质量阈值. q 的取值越高,网络的社区特征就越模糊.由表 3 所示的参数设置可以看出,Network1 与 Network2 拥有同样的网络拓扑结构,但生成的节点属性向量不同(Network3 与 Network4 同理).

表 3 RDyn 网络参数设置

| 网络名称 | n | T | l_{\max} | iter | q | P_{in} | P_{out} | P_{mis} | h |
|----------|-----|-----|------------|------|-----|----------|-----------|-----------|-----|
| Network1 | 400 | 24 | 877 | 80 | 0.2 | 0.5 | 0.05 | 0.1 | 20 |
| Network2 | 400 | 24 | 920 | 80 | 0.2 | 0.5 | 0.1 | 0.2 | 20 |
| Network3 | 400 | 24 | 1106 | 80 | 0.3 | 0.5 | 0.05 | 0.1 | 20 |
| Network4 | 400 | 24 | 1138 | 80 | 0.3 | 0.5 | 0.1 | 0.2 | 20 |

(2)真实动态属性网络

从 DBLP(<https://dblp.org>)抽取论文相关信息来构建真实动态属性网络,这些论文源自 2010 年至 2021 年发表在以下领域的权威学术会议或学术期刊:

(a) 计算机网络: SIGCOMM、MobiCom、NSDI、

INFOCOM、JSAC

(b) 数据库与数据挖掘: SIGKDD、ICDE、ICDM、VLDB、SIGMOD

(c) 计算机视觉与模式识别: CVPR、ICCV、ECCV

收集以上相关论文信息后,首先提取至少有 8 年都发表过论文的作者作为网络节点,然后提取作者间的合著关系作为网络的边,最后每 2 年的网络数据作为一个快照,共有 6 个快照属性网络.每个网络的节点数均统一为 1591,边数依次为 1 468、1 934、2 236、2 347、2 319 和 1 757.对于网络节点属性的构建,则按如下流程:对论文标题进行分词并且去除停用词,过滤掉出现次数小于 10 的词汇,剩余的词汇作为属性集,最终属性数量为 5 052.如果某作者论文标题中的词汇对应属性集中的某个属性,则认为该作者对应的节点包含相应的属性.此外,把作者在相应时间段内发表论文数最多的领域作为其所在网络快照中的真实社区标签.

4.2 对比方法及评价准则

为了对比验证 DAN-NMF 的性能,选择 5 个具有代表性的方法作为基准,包括 DGR-NMF^[15]、Chimera^[19]、TANMF^[21]、RSECD^[24]以及 FaceNet^[28],表 4 归纳了各方法的特性,包括所能够利用的信息种类以及能够完成的任务种类.正如前面指出的是,目前现有方法基本还不能同时完成动态社区发现、社区语义解释及社区演化分析三个任务.

表 4 对比方法的特性

| 方法 | 是否利用拓扑结构信息 | 是否利用属性信息 | 是否利用社区演化平滑约束信息 | 能否检测动态社区 | 能否解释社区语义 | 能否分析社区演化 |
|---------|------------|----------|----------------|----------|----------|----------|
| TANMF | √ | √ | × | √ | √ | × |
| RSECD | √ | √ | × | √ | √ | × |
| DGR-NMF | √ | × | √ | √ | × | × |
| FaceNet | √ | × | √ | √ | × | √ |
| Chimera | √ | √ | √ | √ | √ | × |
| DAN-NMF | √ | √ | √ | √ | √ | √ |

由于所采用的人工合成及真实动态属性网络都具有真实社区标签,因此采用广泛使用的标准化互信息 NMI(Normalized Mutual Information)^[29]作为评价社区发现效果的性能指标,其值越大,则表示社区发现的性能越好.

4.3 实验结果与分析

(1) 参数敏感性分析

DAN-NMF 包含四个权重参数 α 、 β 、 γ 及 η ,需要对这四个参数进行敏感性分析,以确定它们在实验中的最佳取值,这里选择 Dynamic LFR 的 Switch1 和 BirthDeath1 网络作为示例进行分析.对于参数 α 和 β ,首先固定 γ 和 η 的取值为 0,将 α 和 β 的取值范围分别设置为 0 到 50 和 0 到 5,步长分别设置为 5 和 0.5,然后计算 10 个时刻属性网络社区发现结果的平均 NMI,最终结果如图 1(a) 和图 2(a) 所示.可以看出当 α 取值为 20 至 30 并且 β 取值大于等于 2 时,DAN-NMF 都能取得较大的 NMI 值.同时需要指出的是,当 α 取值为 0 并且 β 取值为 0 时,NMI 值较低,这说明融合节点拓扑结构与属性信息,以及考虑社区演化关系能够有效提高社区发现性能.对于参数 γ 和 η ,首先固定 $\alpha=20$ 和 $\beta=2$,将 γ 和 η 的取值范围分别设置为从 0 到 3 和 0 到 5,步长分别设置为 0.3 和 0.5,然后计算 10 个时刻属性网络社区发现的平均 NMI,最终结果如图 1(b) 和图 2(b) 所示.可以看出 γ 取值为 0.3 至 0.9 且 η 取值大于 0.5 时,DAN-NMF 都能取得较大的 NMI 值.当 γ 取值为 0 且 η 取值为 0 时,NMI 值相对较低,这说明考虑社区演化平滑约束信息有利于社区发现性能的提升.

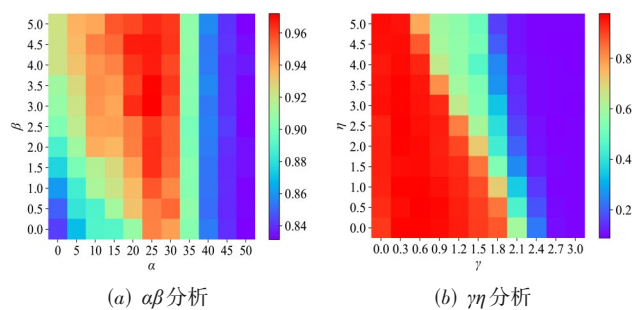


图1 Switch1网络上的参数敏感性分析

在其余动态属性网络进行相同的参数敏感性分析,也具有相似的发现,即 α 取值为 20 至 30, β 取值大于等于 2, γ 取值为 0.3 至 0.9, η 取值大于 0.5 时,DAN-NMF 都能取得较大的 NMI 值.在本文的所有实验中,都统一设置 $\alpha=20$, $\beta=2$, $\gamma=0.3$ 及 $\eta=1$.

(2) 人工合成动态属性网络对比分析

(a) Dynamic LFR 网络. 所有方法在 Dynamic LFR 的 8 个人工合成动态属性网络上分别重复运行 20 次,

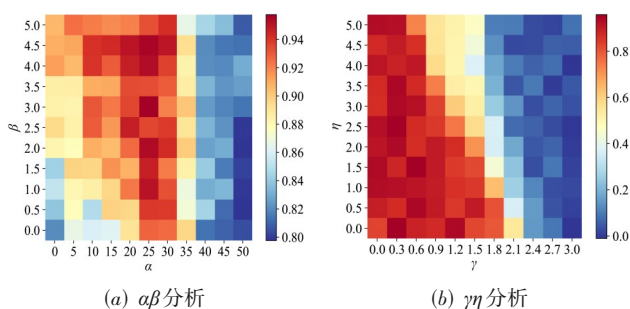


图2 BirthDeath1网络上的参数敏感性分析

各方法的 NMI 取平均值作为评价结果(图 3). 从图 3(a)~(h) 可以看出,对比方法都存在这种问题:在面对某种类型网络时,社区发现性能会有较为明显的下降.比如 RSECD 在 Switch 类型的网络中有着相对较好的性能,但面对 MergeSplit 类型的网络时表现最差.又比如 Chimera,当面对属性信息噪声增加及属性信息与社区结构匹配程度低的网络时,它的性能大多数情况下是最差的.而本文提出的 DAN-NMF 方法无论应用在何种类型的网络,都能在绝大部分时刻获得最好的 NMI,即能够更准确地挖掘出网络潜在的社区结构.这也说明 DAN-NMF 能够充分地利用动态属性网络中蕴含的拓扑结构信息、属性信息以及社区演化平滑约束信息,有效地提升社区发现的准确性及鲁棒性.

DAN-NMF 可以进一步挖掘相邻时刻网络快照的社区演化模式.为更直观地分析社区演化,可以对得到的社区演化矩阵 G_t 进行可视化.作为示例,选取 Switch1 和 Switch2 网络各 2 个社区演化矩阵 G_3 和 G_7 进行分析.由于社区演化矩阵 G_t 的初始化是随机的,所以相同社区前一时刻的编号与后一时刻的编号可能并不相同.为了便于分析,需将相同社区相邻时刻的编号调整为相同,具体调整过程为:对于 $t-1$ 时刻社区编号为 i 的社区,它在时刻 t 中的相应编号通过 $\operatorname{argmax}_l G_{il,t}$ 得到,然后交换 G_t 的第 l 列与 G_t 的第 i 列.其余社区以此类推,从而令相邻时刻的同一社区编号也相同,对调整后的 G_t 进行可视化的结果如图 4 所示.

图 4 各个子图中第 i 行第 j 列格子的颜色深浅表示前一时刻快照的社区 i 转移到下一时刻快照社区 j 的概率.颜色越深,社区 i 转移到社区 j 的概率就越大.从图 4(a) 和图 4(b) 可以看到矩阵对角线基本都呈现深红色,说明 Switch1 网络的社区内大多数节点在下一时刻都会保留到当前社区,社区内只有少部分节点发生了社区转移.而从图 4(c) 和图 4(d) 可以看到矩阵对角线所呈现的红色都相对较浅,说明 Switch2 网络中较多节点发生了社区转移.回顾 Switch1 网络和 Switch2 网络所设置的演化过程:Switch1 每个时刻都会随机挑选 20% 的节点改变其所属社区,剩余 80% 节点保留在原社区,而 Switch2 每个时刻都会随机

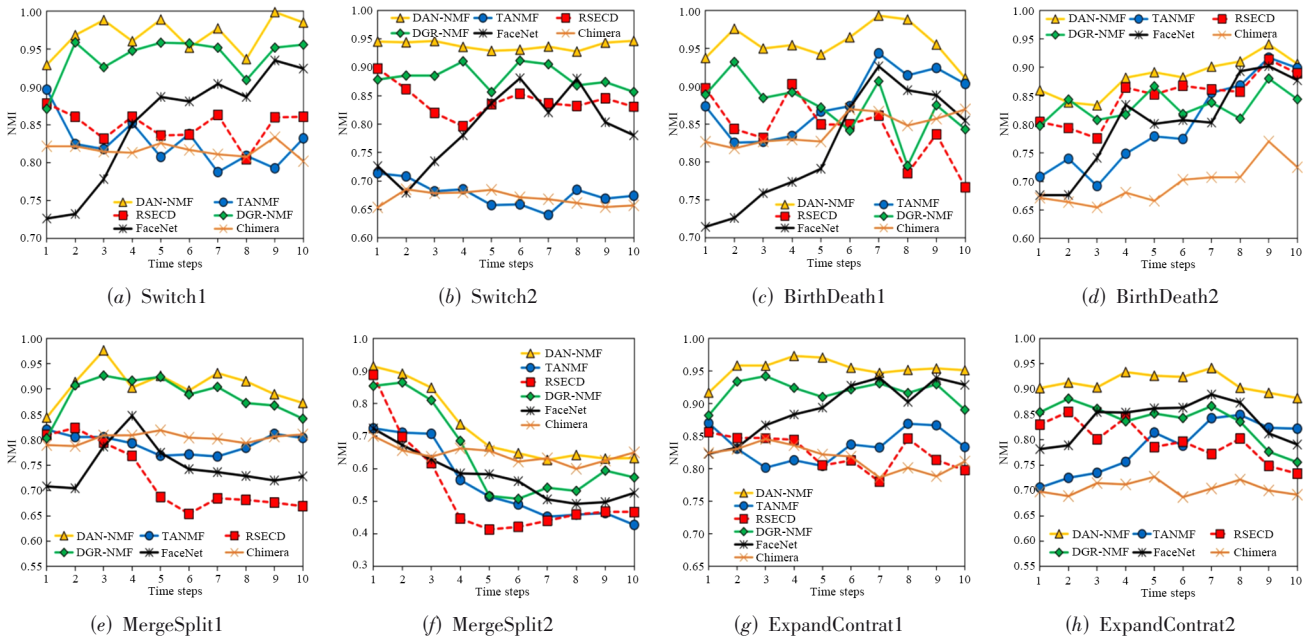


图3 Dynamic LFR 网络的NMI比较

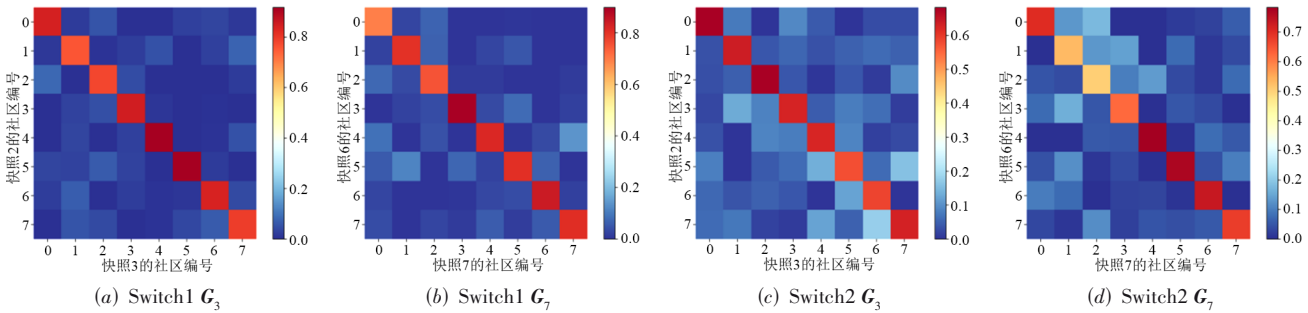


图4 社区演化矩阵 G_t 的可视化

挑选 40% 的节点改变其所属社区, 剩余 60% 节点保留在原社区. 根据前面的观察, 不难发现这两个网络的演化过程与社区演化矩阵所表达的社区间转移概率高度吻合, 这说明 DAN-NMF 的社区演化矩阵 G_t 能准确表示相邻时刻社区间的转移关系.

(b)RDyn 网络. 为进一步验证 DAN-NMF 的普适性, 在另一类人工合成动态属性网络 RDyn 上进行对比实验. 所有方法在 RDyn 网络上重复运行 20 次, 各方法的 NMI 取平均值作为评价结果 (图 5). 通过观察图 5(a) 和图 5(c), 可以看出在 Network1 和 Network3 中, DAN-NMF 只在部分时刻的性能略低于 Chimera, 同时高于其余 4 个对比方法. 进一步地观察图 5(b) 和图 5(d), 可以看出在 Network2 和 Network4 中, DAN-NMF 在绝大多数时刻都取得了最高性能. 这些结果表明 DAN-NMF 具备更强的鲁棒性, 能够适应条件更为苛刻的网络. 需要指出的是, Chimera 的 NMI 变化曲线波动较小, 而其余对比方法 (包括 DAN-NMF) 都存在较明显的波动 (例如在 Net-

work3 和 Network4 上). 这是因为 Chimera 会同时提取所有历史快照网络社区演化平滑约束信息来挖掘当前时刻网络的社区, 从而使得其获得较为平滑的性能, 当然也导致了更高的时间复杂度.

(3) 真实动态属性网络对比分析

为了进一步验证 DAN-NMF 在真实网络同样具备有效性, 在 DBLP 动态属性网络进行实验. 所有方法在 DBLP 网络上重复运行 20 次, 各方法的 NMI 取平均值作为评价结果 (图 6). 通过图 6 可以看出, DAN-NMF 在后四个时刻的网络快照都取得最高的 NMI. 对于 6 个快照网络社区发现结果的平均 NMI, 最优的是 DAN-NMF 的 0.561, 次优的是 Chimera 的 0.523, DAN-NMF 比 Chimera 提高了 7.3%. 以上结果表明 DAN-NMF 在真实网络上也具有较好的社区发现性能.

DAN-NMF 不仅能够发现每个网络快照的社区, 还能给出社区对应的语义解释. 具体地, DAN-NMF 可以计算得到某个社区最关注的主题 $Z=[Z_1, Z_2, \dots, Z_m]$, Z_i 表示词

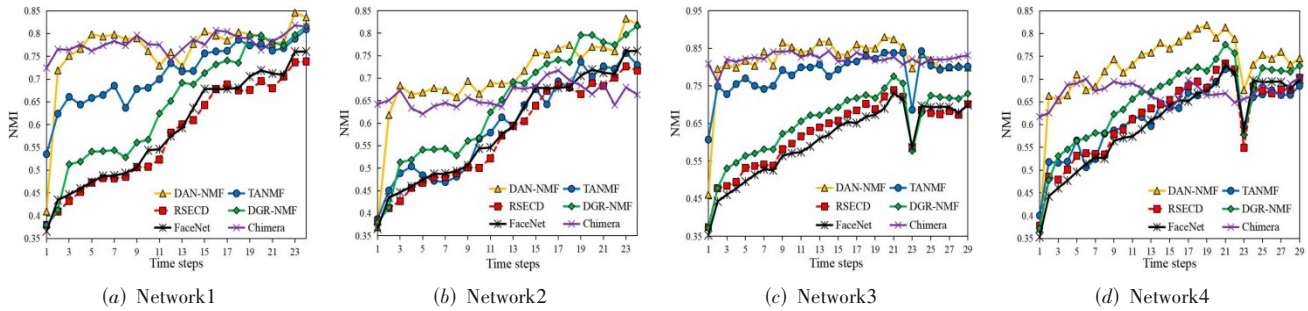


图5 RDyn网络的NMI比较

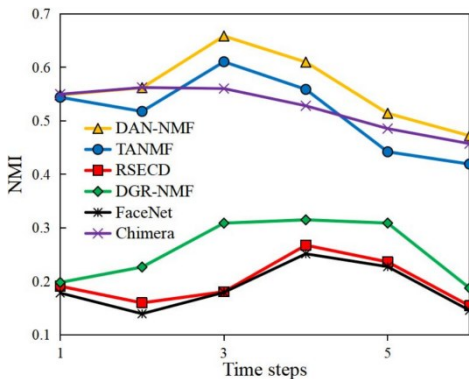


图6 DBLP网络的NMI比较

汇 i 与该主题关联的强度,通过对 Z_i 进行降序排序,并选取排序靠前的词汇作为主题的代表性词汇,最终这些词汇起到解释社区的作用.以2018年至2019年的DBLP网络快照为例,各社区的语义解释如表5所示,词汇后面为其主题隶属强度.通过表5可以看出,社区1的关键词汇包含目标检测(detection)、图像识别(recognition)及语义分割(segmentation),可以推测该社区的研究方向是计算机视觉与模式识别;社区2的关键词汇包含数据库查询(query)和图数据(graph),可以推测该社区的研究方向是数据库与数据挖掘;社区3的关键词汇包含移动通信(mobile)、无线网络(wireless)及云计算(cloud),可以推测该社区的研究方向是计算机网络.

进一步地,结合时间信息可以进一步分析社区语义的变化.假设时刻 t 中社区 i 的代表性词汇集合为 $\text{word}_{i,t}$,通过 $\text{word}_{i,t} - \text{word}_{i,t-1} = \{x | x \in \text{word}_{i,t} \wedge x \notin \text{word}_{i,t-1}\}$ 来获得时刻 t 中社区 i 相对于时刻 $t-1$ 的新增代表性词汇.以表示计算机视觉与模式识别领域的社区为例,计算其在每一网络快照中的新增代表性词汇,结果如表6所示.通过表6可以推断出如下信息:在2014年至2015年这个时间段,深度学习(deep learning)和卷积神经网络(convolutional neural networks)是热门话题;在2018年至2019年这个时间段,生成对抗网络(generative adversarial networks)和注意力机制(attention mechanism)变得流行;在2020年至2021年这个时间段,自监督学习(self-supervised learning)成为新兴热点.

表5 DBLP网络的社区语义解释

| 社区1 | 社区2 | 社区3 |
|-----------------------|--------------------|-----------------------|
| deep (26.74) | efficient (23.89) | networks (66.10) |
| image (24.57) | query (17.52) | mobile (22.92) |
| detection (22.99) | graph (15.01) | wireless (22.83) |
| video (21.20) | search (14.87) | system (15.38) |
| object (20.33) | system (13.26) | edge (13.23) |
| segmentation (19.51) | large (9.94) | distributed (12.36) |
| adversarial (18.13) | extend (9.62) | neural (12.11) |
| neural (17.58) | abstract (9.42) | scheduling (10.86) |
| visual (17.06) | processing (9.38) | communication (10.82) |
| estimation (16.93) | spatial (9.01) | information (10.78) |
| 3D (16.25) | similarity (8.75) | joint (10.41) |
| recognition (16.04) | analytics (8.67) | analysis (10.16) |
| semantic (15.70) | management (8.63) | control (9.68) |
| networks (15.26) | dynamic (8.56) | social (9.36) |
| unsupervised (14.86) | framework (8.30) | adaptive (9.20) |
| feature (13.92) | discovery (8.17) | allocation (9.05) |
| convolutional (13.24) | distributed (7.86) | caching (8.91) |
| attention (12.68) | social (7.84) | resource (8.89) |
| pose (12.56) | fast (7.71) | cloud (8.25) |
| action (11.76) | top-k (7.63) | traffic (8.10) |

表6 计算机视觉与模式识别社区各时间段的新兴热点词汇

| 2012—2013 | 2014—2015 | 2016—2017 | 2018—2019 | 2020—2021 |
|--------------|---------------|--------------|-------------|-----------------|
| mining | convolutional | | generative | self-supervised |
| segmentation | neural | unsupervised | adversarial | |
| information | deep | | attention | |
| | semantic | | | |

综上所述,对于真实世界中的动态属性网络,DAN-NMF不但能够准确地挖掘网络的潜在社区,并对社区做出语义解释,还能够结合时间信息进一步地分析出社区语义的变化情况.

5 结论

由于现有方法难以同时完成动态属性网络语义社区发现及演化分析涉及三个任务:动态社区发现、社区语义解释及社区演化分析,本文为此提出一种基于联合NMF

的方法 DAN-NMF. DAN-NMF 利用联合 NMF 框架集成网络拓扑结构信息、节点属性信息及社区演化平滑约束信息,并通过最大最小化优化策略同时学习社区隶属关系、社区隶属主题及社区演化模式. 在人工合成和真实动态属性网络进行了大量的对比实验,结果表明 DAN-NMF 不仅能高效地发现动态社区,还能挖掘社区演化模式及对社区进行语义解释. 由于 DAN-NMF 方法的时间复杂度近似为 $O(n^2)$,这限制了它应用于大规模动态属性网络的能力. 在后续研究中,将从模型迭代更新算法优化设计及并行化实现这两个方向展开进一步研究,目的在于进一步提高 DAN-NMF 的运行效率.

参考文献

- [1] GARZASE, SCHAEFFER S E. Community detection with the label propagation algorithm: A survey[J]. *Physica A: Statistical Mechanics and its Applications*, 2019, 534: 122058.
- [2] 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.
- [3] MOSCATO V, PICARIELLO A, SPERLÍ G. Community detection based on Game Theory[J]. *Engineering Applications of Artificial Intelligence*, 2019, 85: 773-782.
- [4] JIN D, YU Z Z, JIAO P F, et al. A survey of community detection approaches: From statistical modeling to deep learning[J]. *IEEE Transactions on Knowledge and Data Engineering*, 2023, 35(2): 1149-1170.
- [5] PALLA G, BARABÁSI A L, VICSEK T. Quantifying social group evolution[J]. *Nature*, 2007, 446: 664-667.
- [6] DHOUIOUI Z, AKAICHI J. Tracking dynamic community evolution in social networks[C]//2014 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining (ASONAM 2014). Piscataway: IEEE, 2014: 764-770.
- [7] İLHAN N, ÖÇÜDÜCÜ 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.
- [8] LI W M, ZHOU X K, YANG C, et al. Multi-objective optimization algorithm based on characteristics fusion of dynamic social networks for community discovery[J]. *Information Fusion*, 2022, 79: 110-123.
- [9] AGARWAL P, VERMAR, AGARWAL A, et al. DyPerm: Maximizing permanence for dynamic community detection[M]//Advances in Knowledge Discovery and Data Mining. Cham: Springer International Publishing, 2018: 437-449.
- [10] MENG F R, ZHANG F, ZHU M, et al. Incremental density-based link clustering algorithm for community detection in dynamic networks[J]. *Mathematical Problems in Engineering*, 2016: 1873504.
- [11] 郭昆, 彭胜波, 陈羽中, 等. 基于密度聚类的增量动态社区发现算法[J]. *模式识别与人工智能*, 2018, 31(11): 965-978.
- [12] GUO K, PENG S B, CHEN Y Z, et al. Incremental dynamic community detection algorithm based on density clustering[J]. *Pattern Recognition and Artificial Intelligence*, 2018, 31(11): 965-978. (in Chinese)
- [13] PARK N, ROSSI R, KOH E, et al. CGC: Contrastive graph clustering for community detection and tracking[C]//Proceedings of the ACM Web Conference 2022. New York: ACM, 2022: 1115-1126.
- [14] CHAKRABARTI D, KUMAR R, TOMKINS A. Evolutionary clustering[C]//Proceedings of the 12th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining. New York: ACM, 2006: 554-560.
- [15] 潘剑飞, 曹燕, 董一鸿, 等. 基于 Attention 深度随机森林的社区演化事件预测[J]. *电子学报*, 2019, 47(10): 2050-2060.
- [16] PAN J F, CAO Y, DONG Y H, et al. The community evolution event prediction based on attention deep random forest[J]. *Acta Electronica Sinica*, 2019, 47(10): 2050-2060. (in Chinese)
- [17] 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.
- [18] 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.
- [19] BELLO G A, HARENBERG S, AGRAWAL A, et al. Community detection in dynamic attributed graphs[M]//Advanced Data Mining and Applications. Cham: Springer International Publishing, 2016: 329-344.
- [20] 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.
- [21] APPEL A P, CUNHA R L F, AGGARWAL C C, et al. Temporally evolving community detection and prediction in content-centric networks[M]//Machine Learning and Knowledge Discovery in Databases. Cham: Springer In-

ternational Publishing, 2019: 3-18.

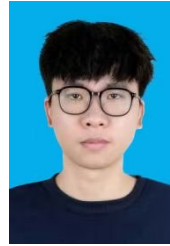
- [20] ZHANG Y L, WU B, NING N W, et al. Dynamic topical community detection in social network: A generative model approach[J]. IEEE Access, 2019, 7: 74528-74541.
- [21] LI D, ZHANG S, MA X. Dynamic module detection in temporal attributed networks of cancers[J]. IEEE/ACM Transactions on Computational Biology and Bioinformatics, 2022, 19(4): 2219-2230.
- [22] LEE D D, SEUNG H S. Learning the parts of objects by non-negative matrix factorization[J]. Nature, 1999, 401: 788-791.
- [23] DING C, HE X F, SIMON H D. On the equivalence of nonnegative matrix factorization and spectral clustering[C]// Proceedings of the 2005 SIAM International Conference on Data Mining. Philadelphia, PA: Society for Industrial and Applied Mathematics, 2005: 606-610.
- [24] 金弟, 刘子扬, 贺瑞芳, 等. 面向带属性复杂网络的鲁棒、强解释性社团发现方法[J]. 计算机学报, 2018, 41(7): 1476-1489.
- JIN D, LIU Z Y, HE R F, et al. A robust and strong explanation community detection method for attributed networks[J]. Chinese Journal of Computers, 2018, 41(7): 1476-1489. (in Chinese)
- [25] CHOUZENOUX E, FEST J B. SABRINA: A stochastic subspace majorization-minimization algorithm[J]. Journal of Optimization Theory and Applications, 2022, 195(3): 919-952.
- [26] 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.
- [27] ROSSETTI G. RDYN: Graph benchmark handling community dynamics[J]. Journal of Complex Networks, 2017, 5(6): 893-912.
- [28] 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.
- [29] CHAKRABORTY T, DALMIA A, MUKHERJEE A, et al. Metrics for community analysis: A survey[J]. ACM Computing Surveys, 2018, 50(4): 1-37.

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



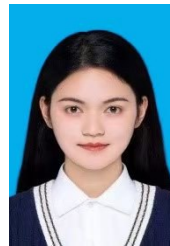
成其伟 男, 1997年8月出生于广东省江门市. 现工作于维沃移动通信有限公司. 主要研究方向为社区发现.

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.

Email: jung@m.scnu.edu.cn



杨佳琦 女, 2000年5月出生于陕西省咸阳市. 现为华南师范大学计算机学院硕士研究生. 主要研究方向为知识图谱.

E-mail: 2022023277@m.scnu.edu.cn



程 颢 男, 2000年2月出生于江西省上饶市. 现为华南师范大学硕士研究生. 主要研究方向是符号图数据挖掘和社交网络分析.

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



汤 庸 男, 1964年5月出生于湖南省张家界市. 现为华南师范大学计算机学院教授, 博士生导师. 主要研究方向为数据智能与云服务、学术社交网络与教育大数据等.

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