

# 控制与决策

Control and Decision

## 一种基于节点嵌入表示学习的社区搜索算法

赵卫绩, 张凤斌, 刘井莲

引用本文:

赵卫绩, 张凤斌, 刘井莲. 一种基于节点嵌入表示学习的社区搜索算法[J]. *控制与决策*, 2021, 36(8): 1970–1976.

在线阅读 View online: <https://doi.org/10.13195/j.kzyjc.2019.1439>

---

## 您可能感兴趣的其他文章

Articles you may be interested in

### 基于种群演化的超参数异步并行搜索

Asynchronous parallel hyperparameter search with population evolution

控制与决策. 2021, 36(8): 1825–1833 <https://doi.org/10.13195/j.kzyjc.2019.1743>

### 基于改进RRT\*FN算法的机器人路径规划

Robot path planning based on improved RRT\*FN algorithm

控制与决策. 2021, 36(8): 1834–1840 <https://doi.org/10.13195/j.kzyjc.2019.1713>

### 面向复杂网络的异常检测研究进展

Research progress of anomaly detection for complex networks

控制与决策. 2021, 36(6): 1293–1310 <https://doi.org/10.13195/j.kzyjc.2020.0055>

### 基于动态行为选择的和声搜索算法

Harmony search algorithm based on dynamic behavior selection

控制与决策. 2021, 36(3): 577–588 <https://doi.org/10.13195/j.kzyjc.2019.0597>

### 一种新的基于标签传播的复杂网络重叠社区识别算法

A novel algorithm for overlapping community detection based on label propagation in complex networks

控制与决策. 2020, 35(11): 2733–2742 <https://doi.org/10.13195/j.kzyjc.2019.0176>

# 一种基于节点嵌入表示学习的社区搜索算法

赵卫绩<sup>1,2</sup>, 张凤斌<sup>1†</sup>, 刘井莲<sup>2,3</sup>

(1. 哈尔滨理工大学 计算机科学与技术学院, 哈尔滨 150080; 2. 绥化学院 信息工程学院, 黑龙江 绥化 152061; 3. 东北大学 计算机科学与工程学院, 沈阳 110169)

**摘要:** 针对已有社区搜索算法采用高维稀疏向量表示节点时间复杂度高的问题, 提出一种基于节点嵌入表示学习的社区搜索算法 CSNERL. 节点嵌入技术能够直接从网络结构中学习节点的低维实值向量表示, 为社区搜索提供了新思路. 首先, 针对已有节点嵌入算法存在较高概率在最亲近邻居间来回游走的问题, 提出基于最亲近邻居但不立即回访随机游走的节点嵌入模型 NECRWNR, 采用 NECRWNR 模型学习节点的特征向量表示; 然后, 采用社区内所有节点的向量均值作为社区的向量表示, 通过选择与当前社区距离最近的节点加入社区的方法实现一种新的社区搜索算法. 在真实网络和模拟网络数据集上分别与相关的社区搜索算法进行实验对比, 结果表明所提出社区搜索算法 CSNERL 具有更高的准确性.

**关键词:** 社区搜索; 节点嵌入; 网络表示学习; 社区发现; 局部社区发现; 随机游走

中图分类号: TP273

文献标志码: A

DOI: 10.13195/j.kzyjc.2019.1439

开放科学(资源服务)标识码(OSID):



**引用格式:** 赵卫绩, 张凤斌, 刘井莲. 一种基于节点嵌入表示学习的社区搜索算法 [J]. 控制与决策, 2021, 36(8): 1970-1976.

## Community search algorithm based on node embedding representation learning

ZHAO Wei-ji<sup>1,2</sup>, ZHANG Feng-bin<sup>1†</sup>, LIU Jing-lian<sup>2,3</sup>

(1. School of Computer Science and Technology, Harbin University of Science and Technology, Harbin 150080, China; 2. School of Information Engineering, Suihua University, Suihua 152061, China; 3. School of Computer Science and Engineering, Northeastern University, Shenyang 110169, China)

**Abstract:** Considering that the existing community search algorithms represent nodes as high-dimensional sparse vectors and have high time complexity, a community search algorithm based on node embedding representation learning (CSNERL) is proposed. Node embedding techniques can learn low-dimensional vectorial representation of nodes from network structure directly, and provide a new solution to community search problems. Firstly, in view of the problem that the existing node embedding algorithm has a high probability to walk back and forth between the closest neighbors, a node embedding model based on closest-neighbor biased random walk with non-immediately revisiting (NECRWNR) is proposed. Based on this model, vectorial representation of nodes is learned and used as feature vectors of nodes in the downstream data mining task. Then, vectorial representation of a community is defined as the average of the vectors for nodes in the community, and a new community search algorithm is designed by choosing those nodes which are nearest to the current community. The proposed algorithm is tested on both real-world and synthetic network datasets with the related community search algorithms. The experimental results show that the CSNERL algorithm is more effective at community search than baselines.

**Keywords:** community search; node embedding; network representation learning; community detection; local community detection; random walk

## 0 引言

社区结构是复杂网络的一个重要属性. 文献[1]将社区结构定义为网络中若干个内部连接紧密的节

点的群组, 节点与所在群组内的节点连接紧密, 与其他群组连接稀疏<sup>[2]</sup>. 挖掘网络中的社区在个性化推荐、信息传播等方面都具有重要意义. 自2002年提

收稿日期: 2019-10-13; 修回日期: 2020-02-18.

基金项目: 国家自然科学基金项目(61172168, 61772122, 61872074); 黑龙江省省属高校基本科研业务费科研项目(YWK10236200141).

†通讯作者. E-mail: zhangfb@hrbust.edu.cn.

出社区发现问题至今,涌现出了基于划分、模块度优化、标签传播等方法的众多社区发现算法,极大地推动了社区发现的发展<sup>[3]</sup>. 这些算法旨在寻找给定网络中的所有社区,但由于时间复杂度高而难以应用于大网络. 同时,许多应用场景也只是对网络中的某一个社区感兴趣,例如已知某研究领域的一名学者,寻找活跃在该领域的其他学者<sup>[4]</sup>;对用户作个性化推荐时,了解其所在社区其他用户的喜好可以提供很多有益信息<sup>[5]</sup>. 不同于传统的社区发现方法,社区搜索<sup>[6-7]</sup>是从给定一个或多个节点出发,寻找包含它们的社区. 相比之下,社区搜索更关注给定节点周围的局部网络结构,能够高效返回用户所关心的社区.

近年来提出了很多社区搜索算法. Clauset<sup>[8]</sup>定义了局部模块度  $R$ ,通过选取使得  $R$  增量最大的外壳节点加入社区,实现了一种含节点集规模  $k$  约束的社区搜索算法. Huang等<sup>[9]</sup>提出基于节点相似度的社区质量度量函数  $tightness$ ,在保证能够增加社区的  $tightness$  值的前提下,从外壳节点集中优先选择与当前社区内节点相似度最高的节点加入社区,计算给定节点所在的社区. 不同于文献[9]采用单层邻居节点集计算节点的相似度, Ma等<sup>[10]</sup>提出采用  $d$ -层邻居节点集计算节点相似度, Liu等<sup>[11]</sup>采用节点嵌入方法将节点映射到低维向量空间,在低维向量空间计算节点间的相似度. 此外, Panagiotakis等<sup>[12]</sup>提出了基于流传播的社区搜索算法.

向量化网络数据是网络数据挖掘首要解决的问题. 传统方法是采用邻接矩阵存储网络,用邻接矩阵的行向量表示节点,但这种高维、稀疏的表示方法导致后续任务的时间复杂度较高<sup>[13]</sup>. 网络表示学习旨在从网络结构中学习节点的低维向量表示,在网络数据和已有数据挖掘模型之间搭起一座桥梁. 受自然语言处理领域词嵌入技术的启发, Perozzi等<sup>[14]</sup>将词嵌入技术引入网络表示学习,提出了第1个节点嵌入模型 DeepWalk,引发了节点嵌入的研究热潮,此后相继涌现出 LINE<sup>[15]</sup>、SDNE<sup>[16]</sup>、node2vec<sup>[17]</sup>、NEMCNB<sup>[11]</sup>等模型.

在此基础上,本文提出一种基于节点嵌入表示学习的社区搜索算法 CSNERL (community search based on node embedding representation Learning). 首先,针对节点嵌入已有工作的不足,提出基于最亲近邻居但不立即回访随机游走的网络嵌入模型 NECRWR (node embedding model based on closest-neighbor biased random walk with non-immediately revisiting),基于该模型,为网络中每个节点学习一个

低维实值向量;然后,采用社区内所有节点的向量均值作为社区的向量表示,通过优先选择与社区距离最近的节点的方法逐渐扩展当前社区,提出基于社区向量表示的社区搜索算法 CSCVR (community search based on community vector representation).

本文主要贡献如下: 1) 提出新的网络嵌入模型 NECRWR. 该模型既考虑了节点间亲近性对随机游走的影响,又避免了在最亲近邻居之间来回游走的问题,可以更好地保留原始网络的性质. 2) 提出的基于节点嵌入表示学习的社区搜索算法 CSNERL 是一个两阶段算法,集网络嵌入模型 NECRWR 与社区搜索算法 CSCVR 于一体,用向量表示社区,将社区外节点与社区的距离建模为两个向量的内积运算,时间复杂度较低. 3) 在真实网络和模拟网络数据集上分别与相关的社区搜索算法和节点嵌入算法进行实验对比,结果表明所提出社区搜索算法 CSNERL 相比 Clauset、GMAC、FlowPro、NEMCNB 等基准算法能够得到更为准确的社区,此外,通过实验验证了 NECRWR 模型有助于提高社区搜索结果的准确性.

## 1 基于节点嵌入表示学习的社区搜索算法

### 1.1 问题定义

用  $G = (V, E)$  表示网络. 其中:  $V$  为网络  $G$  中节点的集合;  $E$  为  $G$  中边的集合,  $\Gamma(v)$  为节点  $v$  的邻居节点的集合. 社区是网络中的一个密集子图,子图内节点连接紧密,子图间节点连接稀疏. 假设对网络  $G$  中的某一个社区  $C$  感兴趣,已知节点  $v \in C$ ,社区搜索问题便是从网络  $G$  中寻找  $k$  个节点,使得这  $k$  个节点中有尽可能多的节点属于社区  $C$ .

基于节点嵌入表示学习的社区搜索问题可以分为两个子问题:节点嵌入表示学习和基于节点向量表示的社区搜索.

**定义1** (节点嵌入表示学习)<sup>[14-15]</sup> 给定网络  $G = (V, E)$ ,节点嵌入的目标是学习一个映射函数  $f: V \rightarrow R^d, d \ll |V|$ ,将网络  $G$  中的任一节点  $u$  映射到  $d$  维向量空间中的一个点  $f(u)$ ,该  $d$  维空间仍保留节点在  $G$  中的性质.

**定义2** (基于节点向量表示的社区搜索) 给定网络  $G = (V, E)$  以及节点在  $d$  维向量空间中的嵌入映射函数  $f$ ,对于目标社区  $C \subset G$ ,已知节点  $v \in C$ ,给定参数  $k$ ,社区搜索的目标是从节点  $v$  出发寻找  $k$  个节点,使得这  $k$  个节点中有尽可能多的节点属于目标社区  $C$ .

### 1.2 CSNERL 算法描述

基于节点嵌入表示学习的社区搜索算法 CSNERL 是一个两阶段算法,集 NECRWR 网络表

示与 CSCVR 社区搜索于一体,算法流程如图1所示.

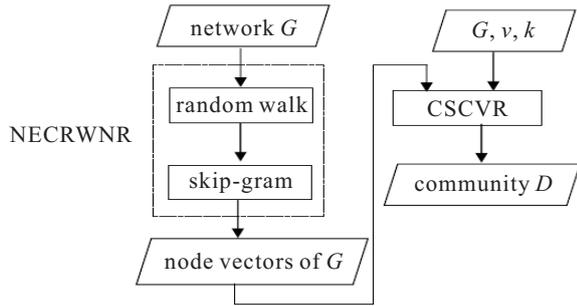


图1 CSNERL算法流程

1.2.1 节点嵌入模型NECRWNR

NECRWNR 与 DeepWalk、node2vec、NEMCNB 同属基于随机游走一类的节点嵌入模型. 首先给出随机游走过程的形式化描述, 然后结合实例分析已有算法的不足.

假设从节点  $v_1$  出发长度为  $l$  的随机游走形成的节点路径表示为  $[v_1, v_2, \dots, v_{i-1}, v_i, v_{i+1}, \dots, v_l]$ , 其中  $v_i$  为该路径对应的节点序列中的第  $i$  个节点. 不失一般性, 将随机游走过程中刚访问过的节点  $v_i$  标记为  $u$ , 则节点  $u$  的邻居节点集  $\Gamma(u)$  中的任一节点  $x$  成为下一个节点  $v_{i+1}$  的概率为

$$P(v_{i+1} = x | v_i = u) = \frac{w_{ux}}{z}. \quad (1)$$

其中:  $w_{ux}$  为边  $(u, x)$  的权重,  $z$  为归一化常数, 计算为

$$z = \sum_{j \in \Gamma(u)} w_{uj}. \quad (2)$$

该类算法之间的主要区别在于计算  $w_{ux}$  的方法不同. DeepWalk 算法将  $w_{ux}$  的值设置为 1; node2vec 算法认为  $w_{ux}$  与节点  $x$  和  $v_{i-1}$  之间的最短路径长度有关; NEMCNB 算法采用 Jaccard Index 度量节点间的相似程度, 认为  $w_{ux}$  与节点  $x$  和  $u$  之间的相似度有关, 在社交网络中人们更倾向于选择与其关系最亲密的邻居节点.

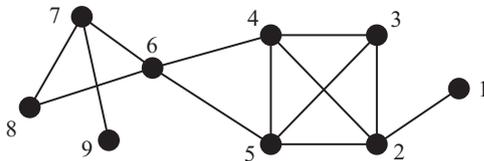


图2 示例网络G

相比 DeepWalk 和 node2vec, NEMCNB 引入节点间相似性的差异指导随机游走, 在社区搜索任务中取得了更好的实验结果, 但在随机游走过程中存在着在最亲近邻居之间来回游走的问题. 结合图2说明该问题, 边的权重  $w$  和游走到邻居节点的概率  $p$  的计算方法参见文献[11], 这里只给出计算结果. 节点4的最亲近邻居是节点5, 从节点4游走到节点5的概率最

大,  $p(5|4) = 0.406$ . 节点5的最亲近邻居是节点4, 从节点5游走到节点4的概率最大,  $p(4|5) = 0.406$ . 因此, 在随机游走过程中从节点4出发, 经过节点5又回到节点4的概率为0.165. 而在基于纯随机的游走方式中, 由于节点4和节点5的度均为4, 从节点4出发到节点5的概率为0.25; 同理, 从节点5出发到节点4的概率也是0.25. 因此, 从节点4出发经节点5再回到节点4的概率仅为0.0625.

通过以上分析可以看出, NEMCNB 算法偏向选择最亲近的邻居节点, 出现在最亲近邻居间来回游走情况的概率较高. 针对此问题, 提出基于最亲近邻居但不立即回访的随机游走 CRWNR (closest-neighbor biased random walk with non-immediately revisiting), 即当新选择的节点  $v_{i+1}$  与  $v_{i-1}$  相同时, 放弃该次选择. 根据各邻居节点的概率分布, 重新从当前节点  $v_i$  的邻居节点集中随机选择一个节点. 算法描述如下所示.

算法1 CRWNR.

输入: 网络  $G = (V, E)$ , 起始节点  $u$ , 行走的节点个数  $l$ ;

输出: 节点路径  $p$ .

step 1: 初始化节点路径  $p = [u]$ .

step 2: 初始化变量  $i = 1$ .

step 3: 由式(1)计算从当前节点  $p[i]$  到其每一个邻居节点的概率.

step 4: 根据  $\Gamma(p[i])$  中每个节点的概率分布, 随机选择一个节点, 标记为  $y$ .

step 5: 如果节点  $y$  与  $p[i-1]$  是同一个节点, 则重复执行 step 4; 否则, 将  $y$  添加到  $p$  中.

step 6:  $i++$ .

step 7: 重复执行  $l-1$  次 step 3 ~ step 6.

step 8: 返回节点路径  $p$ , 算法结束.

step 3 ~ step 5 是本文算法区别于其他随机游走算法的主要之处. step 3 和 step 4 实现了在随机游走过程中偏向最亲近邻居, step 5 实现了随机游走过程中不立即回访.

在 CRWNR 的基础上, 结合 Skip-gram 模型<sup>[14]</sup>, 提出节点嵌入模型 NECRWNR. Skip-gram 模型是自然语言处理中 word2vec 的一种实现方法, 给定一条节点路径  $[v_1, v_2, \dots, v_{i-1}, v_i, v_{i+1}, \dots, v_l]$ , 设定上下文的窗口大小为  $w$ , 则节点  $v_i$  的上下文节点序列  $C(v_i)$  由节点路径中  $v_i$  前后各  $w$  个节点组成, 即  $C(v_i) = [v_{i-w}, \dots, v_{i-1}, v_{i+1}, \dots, v_{i+w}]$ . 将节点看作单词, 将节点路径看作句子, 通过最大化目标函数

$$\max_f \sum_{v_i \in V} \log p(C(v_i)|f(v_i)) \quad (3)$$

学习节点嵌入映射函数  $f$ . 该目标函数是采用当前节点  $v_i$  的向量表示  $f(v_i)$  预测其上下文节点序列  $C(v_i)$  中节点的条件概率. NECRWNR 算法描述如下所示.

#### 算法2 NECRWNR.

输入: 网络  $G = (V, E)$ , 从每一个节点出发的次数  $r$ , 每次行走的节点个数  $l$ , 上下文窗口大小  $w$ , 节点向量的维度  $dn$ ;

输出: 网络  $G$  的节点嵌入映射函数  $f$ .

step 1: 初始化节点路径序列  $ps = [ ]$ .

step 2: 调用 CRWNR( $G, u, l$ ) 完成从  $V$  中的某一个节点  $u$  出发的随机游走, 返回的节点路径标记为  $p$ .

step 3: 将  $p$  加入  $ps$ ;

step 4: 每次从不同的节点出发, 重复执行  $|V|$  次 step 2 和 step 3.

step 5: 重复执行  $r$  次 step 2 ~ step 4.

step 6: 利用  $ps$  中的节点路径构建语料库  $T$ .

step 7: 将  $T, w, dn$  作为输入, 采用 Skip-gram 学习节点嵌入映射函数  $f$ , 返回  $f$ , 算法结束.

step 2 ~ step 5 实现了从网络  $G$  中每一个节点重复进行  $r$  次基于最邻近邻居但不立即回访的随机游走, 得到  $|V| \times r$  条节点路径, 保存于  $ps$  中. 然后, 将节点看作单词, 节点路径看作句子, 把  $ps$  转换为语料库  $T$ , 利用 Skip-gram 算法学习节点嵌入映射函数  $f$ .

#### 1.2.2 基于社区向量表示的社区搜索算法 CSCVR

社区是网络中内部连接紧密的节点的群组. 节点嵌入方法学习的是节点的向量表示, 为了度量节点与社区之间的距离, 首先给出社区的向量表示方法.

**定义3** (社区的向量表示) 在节点向量表示的基础上, 借鉴多维空间中几何中心的思想, 提出采用社区内节点的向量均值作为社区的向量表示. 对于社区  $D$ ,  $D$  内节点个数用  $|D|$  表示, 则社区  $D$  的向量表示  $f(D)$  定义为

$$f(D) = \frac{1}{|D|} \sum_{u \in D} f(u). \quad (4)$$

在社区向量表示的基础上, 给出节点与社区的距离定义.

**定义4** (节点与社区的距离) 在  $f(D)$  的基础上, 节点  $v$  与社区  $D$  的距离  $\text{dis}(D, v)$  定义为

$$\text{dis}(D, v) = 1 - f(D) \cdot f(v), \quad (5)$$

其中运算符  $\cdot$  为向量的内积运算, 其取值范围为  $[0, 1]$ , 描述的是社区  $D$  与节点  $v$  之间的相似程度. 因此,  $\text{dis}(D, v)$  反映节点  $v$  与社区  $D$  之间的距离远近. 通过优先选择与社区连接且距离最近的节点加

入目标社区的方法得到新的社区搜索算法 CSCVR. 算法描述如下所示.

#### 算法3 CSCVR.

输入: 网络  $G = (V, E)$ , 节点嵌入映射函数  $f$ , 初始节点  $v$ , 节点个数  $k$ ;

输出: 最有可能是节点  $v$  所在社区的  $k$  个节点的集合  $D$ .

step 1: 初始化  $D = \{v\}$  以及其外壳节点集  $N = \Gamma(v)$ .

step 2: 初始化社区  $D$  的向量表示  $f(D) = f(v)$ .

step 3: 由式 (5) 计算  $N$  中每一个节点到社区  $D$  的距离, 选择距离  $D$  最近的节点, 标记为  $y$ .

step 4: 将节点  $y$  加入  $D$  中.

step 5: 如果  $D$  内节点个数为  $k$ , 则转 step 9.

step 6: 更新  $D$  的向量表示.

step 7: 更新外壳节点集  $N$ .

step 8: 重复执行 step 3 ~ step 7, 直至  $N$  为空.

step 9: 返回社区  $D$ , 算法结束.

step 3 ~ step 8 通过从社区  $D$  的外壳节点集  $N$  中选择距离  $D$  最近的节点加入社区  $D$  的方法扩展目标社区  $D$ , 直至  $D$  内的节点个数达到  $k$  或者外壳节点集  $N$  为空.

节点  $y$  加入目标社区  $D$  后, 在 step 6 更新社区  $D$  的向量表示. 可用式 (7) 快速计算加入节点  $y$  后的新社区  $D'$  的社区向量表示  $f(D')$ , 有

$$f(D') = \frac{|D|f(D) + f(y)}{|D| + 1}. \quad (6)$$

step 7 的更新外壳节点  $N$ , 可通过增量计算完成  $N$  的更新. 由于节点  $y$  加入到社区  $D$  中导致外壳节点集  $N$  发生变化, 更新后的外壳节点集为  $N' = (N \setminus \{y\}) \cup (\Gamma(y) \setminus D)$ .

算法3的计算主要集中在 step 3 ~ step 8 的  $k - 1$  次循环, 每次循环涉及到  $|N|$  个外壳节点与社区距离的计算. 由于社区和节点都采用  $dn$  维向量表示, 计算外壳节点与社区距离的时间复杂度为  $O(dn)$ . 假设节点的平均度为  $d_{\text{avg}}$ , 那么第  $i$  次循环外壳节点个数的上限为  $i \times d_{\text{avg}}$ , 对应的时间复杂度为  $O(i \times d_{\text{avg}} \times dn)$ .  $k - 1$  次循环的时间复杂度总计为  $O(k^2 \times d_{\text{avg}} \times dn)$ .

## 2 实验与结果分析

### 2.1 实验设定

#### 2.1.1 数据集

分别采用真实网络和模拟网络测试算法的有效性. 真实网络包括: Zachary 空手道俱乐部成员关系网

络 Karate<sup>[18]</sup>、海豚网络 Dolphins<sup>[19]</sup>、美国大学生橄榄球网络 Football<sup>[1]</sup>、美国政治书籍网络 Polbooks<sup>[20]</sup>. 模拟网络由 Lancichinetti 等<sup>[21]</sup>提出的 LFR 模型生成, 该模型的主要参数包括网络中的节点个数  $n$ 、节点平均度  $d_{avg}$ 、节点最大度  $d_{max}$ 、混合参数  $\mu$ .  $\mu$  值越大, 意味着生成网络中的节点有越多的邻居节点在其所在社区之外, 正确识别节点所在社区的难度也就越大. 本文生成 3 组网络数据集, 每组包含 10 个网络. LFR 模型的参数设置情况如表 1 所示.

表 1 LFR 模型的参数值

group name	$n$	$d_{avg}$	$d_{max}$	$\mu$
L10K	10000	10	50	0.05, 0.10, ..., 0.50
L30K	30000	10	50	0.05, 0.10, ..., 0.50
L50K	50000	10	50	0.05, 0.10, ..., 0.50

2.1.2 对比算法及度量指标

为验证 CSNERL 算法的有效性, 与 Clauset<sup>[8]</sup>、GMAC<sup>[10]</sup>、FlowPro<sup>[12]</sup>、NEMCNB<sup>[11]</sup> 进行比较. 此外, 与 DeepWalk<sup>[14]</sup>、node2vec<sup>[15]</sup> 和 NEMCNB<sup>[11]</sup> 比较以验证 NECRWR 模型的有效性.

采用准确率 ( $P$ )、召回率 ( $R$ ) 和  $F$  指标 ( $F_1$ )<sup>[9-11]</sup> 衡量算法的优劣. 评估算法在数据集上的表现时, 对数据集中的每个节点进行一次实验, 取  $n$  ( $n$  为节点个数) 次实验的平均值作为最后结果. 根据  $k$  值设定不

同, 分两种情况度量算法优劣: 1) 同一个起始节点相同  $k$  值时, 采用召回率  $R$  值衡量各种算法的有效性, 召回率  $R$  值越大, 表明算法越好. 2) 同一个起始节点不同  $k$  值时,  $k$  值的大小将导致  $P$  值和  $R$  值产生不同的变化趋势, 为了平衡  $P$  值和  $R$  值, 采用  $F_1$  指标来度量算法的有效性,  $F_1$  指标取值越大, 表明算法越好.

2.2 社区搜索算法比较

在该实验中, 设置  $k$  值为起始节点所在真实社区中节点的个数. 首先在真实网络数据集上进行对比实验, 实验结果如表 2 所示.

表 2 各算法在真实网络数据集上的对比结果

$R$	karate	dolphins	football	polbooks
Clauset	0.8817	0.8841	0.7464	0.7646
GMAC	0.5611	0.6803	0.4808	0.4890
FlowPro	0.7908	0.8889	0.7678	0.7566
NEMCNB	0.5903	0.3754	0.8893	0.4759
CSNERL	<b>0.9739</b>	<b>0.9923</b>	<b>0.9163</b>	<b>0.7949</b>

由表 2 可见, CSNERL 算法在 4 个数据集上均取得了最高的召回率  $R$  值, 表现最好, 其次是 Clauset 和 FlowPro 算法.

在 LFR 模拟网络数据集上进行对比实验. 由于 FlowPro 算法时间复杂度高, L30K、L50K 两组数据集不与该算法进行对比, 比较结果如表 3 所示.

表 3 各算法在 LFR 网络数据集上的对比结果

$R$	Clauset			GMAC			FlowPro			NEMCNB			CSNERL		
	L10K	L30K	L50K	L10K	L30K	L50K	L10K	L30K	L50K	L10K	L30K	L50K	L10K	L30K	L50K
$\mu = 0.05$	0.9140	0.9127	0.9117	0.8806	0.8801	0.8779	0.9695	-	-	0.9708	0.9556	0.9515	<b>0.9991</b>	<b>0.9992</b>	<b>0.9993</b>
$\mu = 0.10$	0.8183	0.8128	0.8119	0.8846	0.8847	0.8866	0.9628	-	-	0.9733	0.9513	0.9480	<b>0.9983</b>	<b>0.9985</b>	<b>0.9990</b>
$\mu = 0.15$	0.7249	0.7180	0.7197	0.8755	0.8751	0.8764	0.9599	-	-	0.9616	0.9444	0.9414	<b>0.9975</b>	<b>0.9983</b>	<b>0.9987</b>
$\mu = 0.20$	0.6399	0.6333	0.6316	0.8680	0.8596	0.8614	0.9521	-	-	0.9527	0.9361	0.9318	<b>0.9958</b>	<b>0.9980</b>	<b>0.9977</b>
$\mu = 0.25$	0.5647	0.5534	0.5526	0.8694	0.8476	0.8411	0.9396	-	-	0.9402	0.9215	0.9163	<b>0.9932</b>	<b>0.9948</b>	<b>0.9941</b>
$\mu = 0.30$	0.4863	0.4866	0.4875	0.8612	0.8273	0.8225	0.9212	-	-	0.9225	0.9029	0.9008	<b>0.9849</b>	<b>0.9835</b>	<b>0.9847</b>
$\mu = 0.35$	0.4206	0.4253	0.4227	0.8276	0.8156	0.8026	0.8879	-	-	0.8901	0.8816	0.8752	<b>0.9649</b>	<b>0.9709</b>	<b>0.9720</b>
$\mu = 0.40$	0.3632	0.3618	0.3598	0.8251	0.7902	0.7827	0.8521	-	-	0.8516	0.8500	0.8433	<b>0.9369</b>	<b>0.9446</b>	<b>0.9430</b>
$\mu = 0.45$	0.3117	0.3054	0.3049	0.7898	0.7683	0.7574	0.7974	-	-	0.8089	0.8064	0.8062	<b>0.8997</b>	<b>0.9015</b>	<b>0.9056</b>
$\mu = 0.50$	0.2471	0.2493	0.2480	0.7305	0.7350	0.7209	0.7198	-	-	0.7467	0.7430	0.7477	<b>0.8347</b>	<b>0.8337</b>	<b>0.8442</b>

由表 3 可以得到如下结论: 1) 随着  $\mu$  值增大, 各种算法呈现下降的趋势, 这也验证了 LFR 网络中  $\mu$  值越大社区发现的难度越高这一先验性质. 2) 算法在  $\mu$  值相同但网络规模不同的数据集上的召回率  $R$  值差别很小, 表明这些算法具有良好的稳定性. 3) 在 3 组 LFR 网络数据集上, CSNERL 算法均取得了最高的  $R$  值, 实验结果最好, 其次是 NEMCNB、

FlowPro、GMAC 算法. 相比之下, 基于局部模块度优化的 Clauset 算法在  $\mu \geq 0.10$  的数据集上表现最差, 且  $\mu$  值越大与其他算法之间的差别也越大.

2.3 节点嵌入方法比较

分别用 DeepWalk、node2vec 和 NEMCNB 学习节点向量, 然后采用基于社区向量表示的社区搜索算法 CSCVR 进行社区搜索, 以此对比各节点嵌入模型在

社区搜索问题中的有效性. LFR网络数据集上的实验结果如表4所示.

由表4可以看出:1)节点嵌入方法结合基于社区向量表示的社区搜索算法 CSCVR 在 LFR 网络上取得了很好的实验结果. 这一方面有 CSCVR 算法的贡献,同时也与节点嵌入算法能够保留原始网络的性质有关系. 2)对比各节点嵌入算法的表现

可以看出,NECRWNR表现最好,其次是NEMCNB、node2vec,相比之下DeepWalk表现最差.这与DeepWalk采用纯随机的游走方式,没有考虑节点间的相似度有关. NECRWNR 优于 NEMCNB 主要源于考虑了节点间的亲近性对随机游走的影响,避免了在最亲近邻居之间来回游走的策略.

表4 节点嵌入方法在LFR网络数据集上的对比结果

R	DeepWalk			node2vec			NEMCNB			NECRWNR		
	L10K	L30K	L50K	L10K	L30K	L50K	L10K	L30K	L50K	L10K	L30K	L50K
$\mu = 0.05$	0.9844	0.9958	0.9959	0.9864	0.9940	0.9938	0.9896	0.9960	0.9961	<b>0.9991</b>	<b>0.9992</b>	<b>0.9993</b>
$\mu = 0.10$	0.9783	0.9919	0.9931	0.9814	0.9928	0.9933	0.9885	0.9945	0.9959	<b>0.9983</b>	<b>0.9985</b>	<b>0.9990</b>
$\mu = 0.15$	0.9712	0.9883	0.9886	0.9801	0.9884	0.9917	0.9866	0.9925	0.9951	<b>0.9975</b>	<b>0.9983</b>	<b>0.9987</b>
$\mu = 0.20$	0.9577	0.9852	0.9861	0.9733	0.9864	0.9882	0.9855	0.9877	0.9889	<b>0.9958</b>	<b>0.9980</b>	<b>0.9977</b>
$\mu = 0.25$	0.9376	0.9811	0.9818	0.9710	0.9849	0.9858	0.9843	0.9858	0.9865	<b>0.9932</b>	<b>0.9948</b>	<b>0.9941</b>
$\mu = 0.30$	0.9078	0.9491	0.9539	0.9540	0.9751	0.9770	0.9758	<b>0.9838</b>	<b>0.9849</b>	<b>0.9849</b>	0.9835	0.9847
$\mu = 0.35$	0.8767	0.9160	0.9244	0.9306	0.9562	0.9577	0.9610	0.9625	0.9637	<b>0.9649</b>	<b>0.9709</b>	<b>0.9720</b>
$\mu = 0.40$	0.8418	0.8846	0.8922	0.9027	0.9316	0.9353	0.9318	0.9357	0.9387	<b>0.9369</b>	<b>0.9446</b>	<b>0.9430</b>
$\mu = 0.45$	0.7985	0.8334	0.8447	0.8589	0.8891	0.8961	0.8935	0.8964	0.8983	<b>0.8997</b>	<b>0.9015</b>	<b>0.9056</b>
$\mu = 0.50$	0.7305	0.7797	0.7877	0.8011	0.8277	0.8372	0.8233	0.8299	0.8376	<b>0.8347</b>	<b>0.8337</b>	<b>0.8442</b>

2.4 运行时间比较

在 L10K 数据集上进行算法运行时间对比实验. 对数据集中每个节点进行一次实验,合计进行 10 万次实验,取平均运行时间进行比较. 所有实验均在 Windows10 操作系统下,由 Python 编码完成,实验所使用的硬件参数是: Intel (R) Core (TM) i7-7 700 CPU @3.60 GHz,内存 DDR4 16G. 比较结果如图3所示.

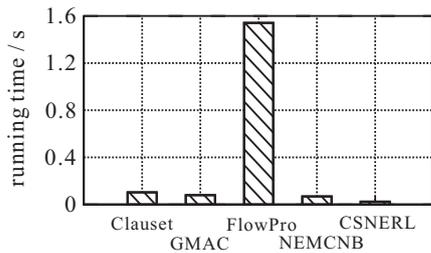


图3 社区搜索算法在L10K上的平均运行时间

由图3可见,除FlowPro进行一次社区搜索平均用时1.54s外,其他算法用时最多的是Clauset的0.10s. 而CSNERL算法平均用时0.02s,远远低于其他4种社区搜索算法.

2.5 参数讨论

讨论参数  $k$  的变化对算法 CSNERL 的影响. 在 L10K ( $\mu = 0.5$ ) 网络上分别使  $k$  取节点所在真实社区节点个数  $n$  的 0.6、0.8、1.0、1.2、1.4 倍进行实验. 由于要对算法在不同  $k$  值的表现进行比较,用  $F_1$  值衡量

算法的好坏. 实验结果如图4所示.

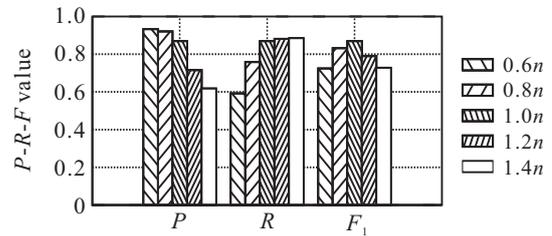


图4 L10K ( $\mu = 0.50$ )数据集上不同  $k$  值的实验结果

由图4可见,在  $k$  逐渐增大的过程中,准确率  $P$  值逐渐减小,而召回率  $R$  值逐渐增大,  $F$  指标  $F_1$  值呈现出先增大后减少的趋势. 较大  $k$  值的返回结果是在较小  $k$  值返回结果的基础上又多返回了一些节点,因此查全率逐渐增大. 由于算法返回的节点个数越来越多,返回结果中位置靠后的节点距离起始节点越来越远,准确率  $P$  值逐渐减小,而  $F_1$  值在  $k$  等于真实社区节点个数时获得最大值.

3 结论

随着互联网和物联网的快速发展,产生了更大规模的网络数据,传统的全局社区发现方法无法有效处理这些大网络. 社区搜索算法由于所需的时间复杂度与所在网络的规模大小没有直接关系,引起了学者们的广泛关注. 受深度学习在网络表示学习方向的启发,本文提出了一种基于节点嵌入表示的

两阶段社区搜索算法CSNERL,并通过实验完成了对CSNERL算法以及NECRWNR模型对社区搜索有效性的验证.此外,观察了参数 $k$ 的变化对CSNERL算法的影响,为参数设置提供了依据.

#### 参考文献(References)

- [1] Girvan M, Newman M E J. Community structure in social and biological networks[J]. Proceedings of the National Academy of Sciences of the United States of America, 2002, 99(12): 7821-7826.
- [2] 邓琨,李文平,陈丽,等.一种新的基于标签传播的复杂网络重叠社区识别算法[J].控制与决策,2020,35(11): 2723-2732.  
(Deng K, Li W P, Chen L, et al. A novel algorithm for overlapping community detection based on label propagation in complex networks[J]. Control and Decision, 2020, 35(11): 2723-2732.)
- [3] Fortunato S, Hric D. Community detection in networks: A user guide[J]. Physics Reports, 2016, 659: 1-44.
- [4] Kloumann I M, Kleinberg J M. Community membership identification from small seed sets[C]. Proceedings of the 20th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining. New York: ACM, 2014: 1366-1375.
- [5] Zhao Y L, Nie L Q, Wang X Y, et al. Personalized recommendations of locally interesting venues to tourists via cross-region community matching[J]. ACM Transactions on Intelligent Systems & Technology, 2014, 5(3): 1-26.
- [6] Sozio M, Gionis A. The community-search problem and how to plan a successful cocktail party[C]. Proceedings of the 16th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining. Washington: ACM, 2010: 939-948.
- [7] Cui W Y, Xiao Y H, Wang H X, et al. Local search of communities in large graphs[C]. Proceedings of the ACM SIGMOD International Conference on Management of Data. Snowbird: ACM, 2014: 991-1002.
- [8] Clauset A. Finding local community structure in networks[J]. Physical Review E, 2005, 72(2): 026132.
- [9] Huang J B, Sun H L, Liu Y G, et al. Towards online multiresolution community detection in large-scale networks[J]. Plos One, 2011, 6(8): e23829.
- [10] Ma L H, Huang H, He Q M, et al. GMAC: A seed-insensitive approach to local community detection[C]. Proceedings of 15th International Conference on Data Warehousing and Knowledge Discovery. Prague: Springer, 2013: 297-308.
- [11] Liu J D, Wang D L, Feng S, et al. Learning distributed representations for community search using node embedding[J]. Frontiers of Computer Science, 2019, 13(2): 437-439.
- [12] Panagiotakis C, Papadakis H, Fragopoulou P. Local community detection via flow propagation[C]. Proceedings of IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining. Paris: ACM, 2015: 81-88.
- [13] Cui P, Wang X, Pei J, et al. A survey on network embedding[J]. IEEE Transactions on Knowledge and Data Engineering, 2019, 31(5): 833-852.
- [14] Perozzi B, Al-Rfou R, Skiena S. DeepWalk: Online learning of social representations[C]. Proceedings of the 20th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining. New York: ACM, 2014: 701-710.
- [15] Tang J, Qu M, Wang M, et al. LINE: Large-scale information network embedding[C]. Proceedings of the 24th International Conference on World Wide Web. Florence: ACM, 2015: 1067-1077.
- [16] Wang D X, Cui P, Zhu W W. Structural deep network embedding[C]. Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining. San Francisco: ACM, 2016: 1225-1234.
- [17] Grover A, Leskovec J. Node2vec: Scalable feature learning for networks[C]. Proceedings of the 22th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining. San Francisco: ACM, 2016: 855-864.
- [18] Zachary W W. An information flow model for conflict and fission in small groups[J]. Journal of Anthropological Research, 1977, 33(4): 452-473.
- [19] Lusseau D, Schneider K, Boisseau O J, et al. The bottlenose dolphin community of doubtful sound features a large proportion of long-lasting associations—Can geographic isolation explain this unique trait?[J]. Behavioral Ecology and Sociobiology, 2003, 54(4): 396-405.
- [20] Newman M E J. Modularity and community structure in networks[J]. Proceedings of the National Academy of Sciences, 2006, 103(23): 8577-8582.
- [21] Lancichinetti A, Fortunato S, Radicchi F. Benchmark graphs for testing community detection algorithms[J]. Physical Review E, 2008, 78(4): 046110.

#### 作者简介

赵卫绩(1980—),男,副教授,博士生,从事社区发现、数据挖掘等研究, E-mail: sdzhaoweiji@163.com;

张凤斌(1965—),男,教授,博士生导师,从事网络安全、入侵检测技术等研究, E-mail: zhangfb@hrbust.edu.cn;

刘井莲(1980—),女,副教授,博士生,从事社会网络分析、社区发现等研究, E-mail: datamining@163.com.

(责任编辑: 郑晓蕾)