# 基于拓扑结构表示学习的大规模无监督图对齐方法研究

王晨旭11.20 周俊铭1 姜佩京1)

<sup>1)</sup>(西安交通大学软件学院 西安 710049)

2)(西安交通大学智能网络与网络安全教育部重点实验室 西安 710049)

**摘 要** 图数据因其较强的复杂关系表征能力受到广泛关注,在社交网络、学术合作、道路交通、生物信息等多个领域具有重要应用.图对齐技术旨在找出不同图中属于同一实体的节点对,在多个领域具有重要的应用价值,例如,对不同社交网络中属于同一个用户的账号进行关联可以为推荐系统提供更丰富的用户行为画像,对不同生物组织的蛋白质网络进行对齐能够辅助研究人员分析蛋白质的特性和机能.然而,在缺乏人工标注信息的情况下仅使用图的拓扑结构信息实现无监督图对齐一直是图数据挖掘面临的重要难题之一,特别是在大规模图对齐任务中,存在初始种子节点发现难和计算效率低下的问题.针对以上问题,本文提出了一种基于拓扑结构表示学习的大规模无监督图对齐框架.首先认待匹配图中选取典型子图作为种子节点候选集,利用局部拓扑结构信息计算得到高可靠的种子节点匹配对,然后利用所得种子节点将待匹配图进行融合,并提出一种高效的无监督表示学习算法将融合图映射到统一的向量空间中,最后利用学习得到的节点向量实现大规模图对齐.与已有方法相比,本文所提方法在大规模图对齐任务中用时最短,对齐结果准确率最高,且算法性能受图结构的差异性影响最小.

关键词 图对齐;大规模图;无监督表示学习 中图法分类号 TP301 DOI号 10.143 P.J.1016.2023.01350

# Large-Scale Graph Alignment Based on Topological Structure Representation Learning

WANG Chen-Xu<sup>1),2)</sup> ZHOU Jun-Ming<sup>1)</sup> JIANG Pei-Jing<sup>1)</sup>

<sup>1)</sup> (School of Software Engineering, Xi'an Jiaotong University, Xi'an 710049) <sup>2)</sup> (Ministry of Education Key Lab of Intelligent Network and Network Security, Xiran Jiaotong University, Xi'an 710049)

**Abstract** Graph data has attracted lots of attention due to its strong ability of representing complex relationships. It has been widely used in many fields, such a social networks, academic cooperation, road traffic, and biological information. Graph alignment aims to find node pairs belonging to the same entity in different graphs, which have valuable applications in many fields. For example, associating accounts belonging to the same user in different social networks can provide richer user behavior profiles for recommender systems, and aligning protein networks of different biological tissues can assist researchers in analyzing the characteristics and functions of proteins. However, unsupervised graph alignment using the topological information of graphs has always been one of the important problems faced by graph data mining in the absence of manual annotation information. There are difficulties in finding initial seed nodes and low computational efficiency, especially for large-scale graph alignment tasks. To solve these problems, this paper proposes a large-scale unsupervised graph alignment framework based on topological structure representation learning. Firstly, a typical subgraph is selected from each of the graphs as a candidate

收稿日期:2022-07-18;在线发布日期:2023-01-09.本课题得到国家自然科学基金(62272379)、陕西省自然科学基金(2021JM-018)、国家 重点研发计划(2021YFB1715600)、中央高校基本科研业务费专项资金(1191320006)资助.**王晨旭**(通信作者),博士,副教授,中国计算机 学会(CCF)会员,主要研究方向为图大数据挖掘、人工智能、网络安全等.E-mail: cxwang@mail.xjtu.edu.cn.**周俊铭**,硕士,主要研究方 向为图匹配对齐、网络表示学习等.**姜佩京**,硕士研究生,主要研究方向为图匹配对齐、网络表示学习等.

set of seed nodes. The local topological information is used to retrieve a set of highly reliable seed node pairs. Then, we use the seed nodes to fuse the matching graphs, and propose an efficient unsupervised representation learning algorithm to map the fused graph into a unified vector space. Finally, large-scale graph alignment is realized based on the learned node vectors. Compared with existing methods, the proposed approach uses the least time in large-scale graph alignment tasks and achieves the best performance of alignment accuracy. Moreover, the structural differences of graphs have limited impacts on the performance of the proposed method.

Keywords graph alignment; large-scale graphs; unsupervised representation learning

# 1 引 言

图数据结构作为对复杂关系进行建模的重要形 式之一,在许多领域有着重要的应用,如在线社交媒 体中的社交网络[1]、科学研究中的合作者网络[2]、生 物信息领域的蛋白质网络<sup>[3]</sup>等. 图数据挖掘能够获 取很多有价值的信息,如节点分类可以推断科研人 员所属的研究领域<sup>[4]</sup>、节点聚类可以分析社交网络 中的兴趣社区<sup>[5]</sup>、链接预测能够为用户推荐更多的 好友关系[6]、节点重要性评估能够分析蛋白质在生、 物功能中的重要性「了及社交网络中用户的影响力」 等.近年来,旨在识别不同图中的节点是否属于同一 🗸 实体的图对齐任务受到越来越多的关注<sup>[9]</sup>,例如,图 对齐能够识别不同社交网络中的账号是否来自同一 用户,帮助推荐系统构建更为完善的用户画像,提高 推荐的准确度[10];图对齐能够发现不同生物蛋白质 网络中相同或相似的蛋白质关联关系,为生物科学 研究提供重要的参考价值[11].

图表示学习又名图嵌入<sup>[4]</sup>,旨在将图中的节点 表示为低维、实值、稠密的向量形式,得到的向量表 示在空间中的距离反映了节点在图中的亲和度.这 些向量可作为特征输入到传统机器学习算法模型 中,进而完成一些常见的下游任务,如节点分类、链 接预测及社区发现等.近年来,利用图表示学习来完 成图对齐任务逐渐成为研究的热点.

由于图对齐任务的重要应用价值,已有不少图对 齐算法被提出,目前图对齐算法大致可以分为两种, 第一种是传统的离散型图对齐算法:如 BigAlign<sup>[12]</sup> 认为图对齐是对源图和目标图的邻接矩阵的行和列 进行重新排序,并将图对齐问题转化为置换矩阵的 最优求解问题;IsoRank<sup>[13]</sup>利用节点邻域拓扑结构 的相似性实现蛋白质网络的对齐;FINAL<sup>[14]</sup>提出结 构相似度、节点特征相似度和边特征相似度三个衡

量标准来解决图对齐问题. 第二类为基于图表示学 习的对齐方法:如 REGAL<sup>[15]</sup>首先通过计算邻居节 点度的相似性来学习节点的向量表示,并利用节点 向量的相似性来完成图对齐任务; CrossMNA<sup>[16]</sup> 通过图间向量的方式实现多个图的对齐任务,该方 法通过图间信息不断优化2类节点向量,并通过图 间向量的相似性实现节点对齐;GAlign<sup>[17]</sup>采用一种 基于多层向量表示的无监督图对齐框架,该方法通 过图卷积网络(GCN)学习节点的向量表示,提出一 种数据增强方法确保模型满足图的一致性约束; WAlign<sup>[18]</sup>使用轻量级的 GCN 结合 Wasserstein 距 离鉴别器,确保在训练的过程中源图和目标图的节 点具有相似的分布概率,从而找到最佳的节点对应 关系.这些算法在取得成功的同时,也面临着重大的 挑战.首先,随着当前图数据的规模不断扩大,已知 的图对齐算法计算效率往往较低,节点之间复杂的 关联关系又导致了算法难以并行化,导致很多算法 无法应用到大规模图数据的对齐任务上.其次,传统 的图对齐算法往往针对单一领域的图数据而设计, 缺乏通用性.如 IsoRank 算法专门为蛋白质网络对 齐而设计[13],将其直接应用到社交网络的对齐任务 时往往出现性能较低、匹配结果不可信的问题.

目前图表示学习在大规模图上进行无监督匹配 对齐主要存在两个难点:一是初始种子节点发现难, 由于不存在先验的种子匹配节点,直接对全图进行 匹配会导致较高的错误率,因此需要先找到一部分 高可靠的匹配节点对作为种子.但如何找到有效的 高可靠种子节点以及如何实现初始匹配是研究的难 点.二是将种子节点应用到全图匹配时,当前大部分 的图表示学习算法均存在计算效率与表示学习有效 性之间的矛盾.部分图表示学习方法仅使用图的局 部结构信息,只保留了节点的1阶或2阶邻居相似 性信息,尽管这种做法时间开销较小,但是难以对全 图进行有效的表示学习;部分方法使用了图的高阶 邻居信息,能够学习到距离较远节点之间的相似性 关系,相较于仅使用图的局部结构信息而言,这种做 法可以提升图对齐的准确率,但在捕获全局结构信 息时往往存在计算效率低下、时间开销大的问题.

当前已有的基于图表示学习的对齐算法往往分 别学习 2 个待匹配图的节点向量表示<sup>[10,13]</sup>,然后利 用学习得到的表示向量计算节点间的相似性,相似 度高的节点对会被判定为匹配节点.然而,当前基于 拓扑结构的图表示学习算法往往需要对节点表示进 行随机的初始化,导致学习得到的待匹配图的节点 表示向量未分布在同一向量空间<sup>[10,19]</sup>,需要将二者 映射到同一空间中,特别是在无监督的条件下,往往 会导致较低的匹配准确率.

为了解决以上问题,本文提出了一种基于图表 示学习的两阶段图对齐框架 LGA (Large-scale Graph Alignment),首先采用高可靠的无监督图匹 配算法对两个子图进行匹配得到种子节点,然后利 用种子节点对2个大图进行融合,然后对融合图进 行无监督表示学习得到节点的表示向量,如此可以 直接将待匹配图映射到统一的向量空间中,避免了 由于表示空间的转换带来的匹配误差.具体步骤如 图1所示.第一阶段为高可靠种子节点发现,先从要 匹配的大图中抽取它们的典型子图,利用少量有代 表性的节点构成典型子图,然后采用高可靠的无监 督图匹配算法对两个子图进行匹配得到种子节点. 第二阶段则利用种子节点对2个大图进行匹配,先 是根据种子节点对2个大图进行融合,将二者合为 一个更大规模的图,然后对融合图进行无监督表示 学习得到表示向量,最后利用表示向量计算节点之 间的相似性进行匹配对齐.由于融合图的节点数量 非常庞大,目前已有的图表示学习算法无法快速高 效地表示学习.针对此问题,本文提出一种基于稀疏 矩阵分解和个性化随机游走相结合的高效图表示学 习方法,在保证计算效率的同时最大限度地捕获融 合图的全局结构信息,为了验证算法的有效性,本文 在 6 个不同类型和不同规模的数据集上进行实验验 证,并与多个已知最先进的无监督图对齐算法进行 对比,实验结果表明,本文所提方法能够快速高效地 实现大规模图的匹配对齐,取得了最高的匹配准确 率,且在两个匹配图结构差异较大的情况下依然具 有较高的匹配准确率.本文工作的主要贡献如下:

(1)提出一种基于种子节点融合的大规模图对 齐算法,与以往利用种子节点监督图表示学习的范 式不同,本文所提方法直接在融合图上进行无监督 的图表示学习,直接将待匹配的图数据映射到统一 的向量空间中,避免了由于表示空间转换导致的匹



图 1 LGA 大规模图对齐框架示意图

配准确率急剧下降的问题.

(2)提出一种基于稀疏矩阵分解和个性化随机 游走相结合的图表示学习方法,实现了对大规模图 的高效表示学习,得到的表示向量能够捕获图的全 局结构信息,并能够很好地适用于图对齐任务.

(3) 在多种类型和不同规模的数据集上进行了 广泛的实验验证,结果表明本文所提方法能够应用 于大规模图数据的高效对齐,取得了远高于其他方 法的匹配准确率和计算效率.本文算法代码已开源: https://github.com/zhoujunming2019/LGA-master.

## 2 相关工作

### 2.1 传统图对齐算法

图对齐在多个领域都有着重要应用,如在生物 领域可以研究蛋白质之间的关系<sup>[11]</sup>,在社交网络领 域<sup>[12]</sup>,跨平台的社交网络用户对孟可以用于用户画 像、用户兴趣挖掘;在知识图谱领域可以实现知识图 谱的补全和扩展<sup>[9]</sup>.传统图对齐算法大多**有读**计算 对齐矩阵 *S*,通常以邻接矩阵的形式输入源图和目 标图,在对齐过程中通过优化损失函数不断地更新 对齐矩阵<sup>[7]</sup>.

BigAlign 将图对齐问题建模为目标图邻接矩 阵的行与列的重新排列问题<sup>[12]</sup>,通过优化损失函数 找到置换矩阵 H 使得排列后的目标图邻接矩阵 接近于源图的邻接矩阵.为了解决求解过程中存在 的 NP(Non-deterministic Polynomial-time)难问题, BigAlign 使用更为宽松的约束,并在优化的过程中 使用迭代优化提高计算效率;IsoRank 认为如果两 个图中节点的邻域相似,那么它们将大概率是匹配 的<sup>[13]</sup>.因此,IsoRank 认为不同蛋白质网络中的节 点在它们各自的序列和邻域拓扑结构相似的情况下 将获得较高的匹配率,基于该思路提出了一种专门 针对蛋白质网络的节点对齐算法.然而传统图对齐 方法除了对结构噪声较为敏感外,其计算效率也不 尽如人意.

### 2.2 基于图表示学习的图对齐方法

随着图表示学习技术的发展,基于节点表示的 图对齐算法逐渐成为研究的热点.IONE<sup>[8]</sup>在研究社 交网络用户对齐问题中提出将每个用户的粉丝 (follower-ship)/关注者(followee-ship)建模为输 入/输出的上下文表示,以保持拥有相似粉丝/关注 者的用户的表示向量在嵌入空间中距离较近.为了 进一步解决不同图中节点所在嵌入空间不一致的问 题,IONE 选取具有较大匹配概率的节点作为锚节 点,使其对相邻的节点在空间中提供的上下文嵌入 表示权值更高,最后采用随机梯度下降和负采样算 法解决扩展性问题.REGAL<sup>[11]</sup>通过对节点一阶近 邻和二阶近邻的度信息进行加权计算获取节点初始 的表示向量,根据初始表示向量的相似性计算得到 节点的相似度矩阵,使不同图中具有相似结构的节 点具有更高的相似度,最终对相似度矩阵进行分解 从而得到节点的低维向量表示.此方法仅考虑了图 的局部拓扑信息,通常情况下对结构噪声的较为 敏感.

DeepLink<sup>[20]</sup>基于随机游走策略生成结点序列 后使用 Skip-Gram 模型对目标图和源图分别训练 得到两个不同的嵌入向量空间,在解决特征空间不 一致性的问题时引入深度神经网络,利用种子节点 训练得到两个特征空间的非线性变换关系.为了提 高图对齐的准确率,DeepLink 在设计损失函数时考 虑了对偶学习,在有监督的方式下使损失函数最小 化来学习源图和目标图的向量表示,需要大量的人 工标注信息,往往难以获得. Chu 等人<sup>[16]</sup>针对多个 图对齐问题提出了 CrossMNA 算法,该模型利用图 之间的交叉信息不断优化节点的向量表示,将图 间向量用于图对齐任务,图内向量用于其他下游图 分析任务:实验结果表明 CrossMNA 能够取得较好 的图对齐效果,并在较少内存使用的情况下实现更 好的连接预测任务. Trung 等人<sup>[17]</sup>提出了 GAlign 模型,该模型利用图卷积网络(Graph Convolutional Networks, GCN)的多层嵌入并结合拓扑一致性原 则实现无监督图对齐在务,在训练过程中,通过比对 GCN 各层之间的嵌入向量来获取锚节点,通过锚节 点实现 GCN 的加权传播. 此外,该模型引入基于扰 动的数据增强方法和自适应的损失函数来增强模型 的鲁棒性, Gao 等人<sup>[18]</sup>提出了 WAlign 模型,该模型 能够捕获图的局部和全局结构信息,并解决了传统 GCN 由于卷积层数的增加而导致的过度平滑问题. 该模型采用一种轻量级的基于 GCN 的图对齐框 架,在优化过程中使用 Wasserstein 距离鉴别器,确 保在训练的过程中源图和目标图的节点具有相似的 概率分布,并引入新的损失函数来避免训练过程中 得到平凡解,找到最优的节点对应关系.

Derr等人<sup>[21]</sup>提出了一种基于深度对抗学习的无 监督图对齐方法 DANA(Deep Adversarial Network Alignment),该算法通过对齐节点的表示向量分布 来发现复杂的节点映射关系,并在此基础上实现快 速的最近邻节点对齐.此外,该模型还具备一定模型 选择能力.Nguyen等人<sup>[22]</sup>提出了 NAWAL 模型, 该模型将源图和目标图分别投影到各自的向量空间 中,并使用一种基于对抗学习的无监督对齐方法来 优化源图和目标图之间的映射函数,从而找到最佳 的匹配映射.实验表明 NAWAL 能够在真实数据集 上取得不错的对齐效果,并且可以与现有的有监督 方法相结合,利用先验信息更好地指导对齐.Zheng 等人<sup>[23]</sup>提出了 JORA 模型来解决社交网络间用户 身份的弱监督对齐问题,该模型同时优化了表示学 习和图对齐模型,同时保留了图内和图间节点的相 似性信息,并采用注意力机制自适应地学习节点的 相似性,极大地减少了模型对标签数量的依赖.

与以上方法相比,本文提出一种新的基于融合 图表示学习的大规模无监督图对齐框架.首先采用 典型子图获取高可靠的种子匹配节点对,然后利用 种子节点对目标图和源图进行融合,并提出了一种 高效的图表示学习算法直接将融合图映射到统一的 向量空间,避免了由于特征空间转换带来的误差.且 不需要对深度神经网络模型进行训练,在保证计算 效率的同时尽可能地捕捉图的结构信息,大大提高 了准确性和鲁棒性.

## 3 高可靠种子节点挖掘

设待匹配图为  $G_1(V_1, E_1)$  和  $G_2(V_2, E_2)$ ,其中  $V_1$ 和 $V_2$ 分别为两个图的节点集合, $E_1$ 和 $E_2$ 为边的 集合,为了表达上更清晰方便,本文在后续章节均采 用符号 u 代表图  $G_1$  中的节点,采用 v 表示图  $G_2$  中的 节点.基于拓扑结构的无监督图对齐可以定义为:在 两个节点子集  $V_1 \subset V_1$  和  $V_2 \subset V_2$  找到一个映射  $\mathcal{M} = \{(u,v) \mid u \in V_1^s, v \in V_2^s\}, 使得由 V_1^s 和 V_2^s$ 构成 的两个子图之间具有最相似的拓扑连接结构.在缺 乏先验信息和节点属性信息的情况下,对  $G_1$ 和  $G_2$ 进行无监督对齐主要依赖于二者之间的拓扑结构信 息,在图规模较大的情况下直接进行全图匹配往往 计算效率较低且会存在较高的错误率,为此,本文首 先在图中选取典型子图作为候选种子节点集,对候 选种子节点进行对齐,得到源图 G<sub>1</sub>与目标图 G<sub>2</sub>中 高可靠的匹配节点对作为种子节点,为下一步的匹 配计算奠定基础.

#### 3.1 典型子图抽取

由真实世界数据构建的图结构,其节点的度往 往遵循幂律分布<sup>[4]</sup>,即只有少部分节点拥有大量的 邻居节点而大部分节点只有较少的连接.图2为对 本文后续实验中涉及到的6个公开数据集节点度的 分布情况,由图可知,所有数据集均呈现出少部分 节点拥有大量连接的情况,这部分节点包含了更加 丰富的结构信息,首先对这些节点进行匹配对齐所 得结果的准确性和可靠性会更高.因此,本文首先分 别从待匹配图中选取度最大的前 K 个节点构成典 型子图,然后对典型子图进行无监督图对齐得到种 子节点.对待匹配图  $G_1(V_1, E_1)$ 来说,选取度数最大 的前 *K* 个节点构成集合  $V_{1'} = \{u_1, u_2, \cdots, u_K\}, 则$  $V_{1'}$ 为 $V_1$ 的一个子集,设由 $V_{1'}$ 中的元素构成的子图 为 $G_{1'}(V_{1'}, E_{1'})$ ,其中 $E_{1'}$ 为原始边集合 $E_1$ 的子集, 且仅保留 V<sub>1</sub>′中节点之间的边,如此得到的子图  $G_{1'}(V_{1'}, E_{1'})$ 即为 $G_1(V_1, E_1)$ 的典型子集,由于 $V_{1'}$ 中的节点在原始图中度最大,因此典型子图  $G_{1'}(V_{1'})$ ,  $E_{1'}$ )保留了原始图的大部分拓扑结构信息.同样地, 设 $G_2(V_2, E_2)$ 的典型子图为 $G_{2'}(V_{2'}, E_{2'})$ .本节所要 挖掘的高可靠种子节点则是由 V1/和 V2/的子集构成 的节点对的集合,该集合也是M的一个子集.

为了使 K 能够与原始图规模相适应,设 K =  $\beta\sqrt[3]{n}$ ,其中 n 代表图中节点数量, $\beta$ 为可调参数.种子 节点数 K 的选择主要考虑了算法的普适性,即无论 图规模的大小如何,算法均能够高效准确地选择合 适的种子节点数.如果按照百分比选取种子节点,对 于大规模图来说会造成种子匹配准确率下降以及预 匹配效率低了的问题,而对小型图来说则会导致种 子节点过少难以发挥预匹配的作用.本文提出的计 算公式既不会在大规模网络上选择过多的节点进行 预匹配,也不会在小型图上选择过少的节点来构建 典型子图,从而达到算法适配的效果.



图 2 不同图中节点度的累计概率分布

#### 3.2 典型子图匹配

本文采用 DeepMatching 算法<sup>[19]</sup>对典型子图进

行无监督对齐以得到高可靠种子节点,首先采用 DeepWalk<sup>[24]</sup>图表示学习算法得到两个典型子图的节 点向量表示,然后采用一致性点漂移算法(Coherent Point Drift,CPD)<sup>[25]</sup>得到待匹配节点之间的匹配概 率矩阵,最后采用最大权匹配算法得到高可靠的种 子节点.

3.2.1 典型子图表示学习

设抽取得到的典型子图为 $G_{1'}(V_{1'}, E_{1'})$ 和 $G_{2'}(V_{2'}, E_{2'})$ ,首先使用 DeepWalk<sup>[24]</sup>算法对两个子图进行编码得到节点的向量表示,DeepWalk 算法包括 2 个步骤.第一步采用随机游走算法生成节点的上下文邻接节点信息,随机游走过程中每次从所在节点的邻居中随机选择一个节点作为下一跳,不断循环直到满足一定的路径长度后停止游走.第二步采用word2vec 技术中的 Skip-Gram 模型来学习节点的向量表示,在此过程中将图中的节点视为 word2vec 语言模型中的单词,将随机游走所得的路径视为自然语言中的句子,该模型通过最大化节点的上下文概率来学习节点的向量表示.设由 DeepWark 得到的两个典型子图的节点向量表示分别为  $R^*$   $H^*$ .

得到节点的向量表示后,本文使用一致性点 移算法(Coherent Point Drift, CPD)<sup>[25]</sup>计算种子节 点之间一致性的概率矩阵.一致性点漂移算法将两 个点集 **R**<sup>1</sup>和 **R**<sup>2</sup>中节点之间的对应关系建模为关联 概率问题,使用概率值描述两个点集中节点之间存 在对应关系的确定性程度,即节点之间的匹配概率. CPD使用混合高斯模型(Gaussian Mixture Model, GMM)描述点集 **R**<sup>1</sup>和 **R**<sup>2</sup>中所有节点对之间的确定 性关系,然后对高斯混合模型的参数进行求解获得 节点之间的匹配关系,通过最大似然法拟合高斯混 合模型的质心,使概率密度函数的期望最大化,求解 混合高斯模型中参数可得到节点间的匹配概率. GMM 模型采用 EM 算法来计算两个图中种子节点 之间的匹配概率矩阵 **P**,在期望求解阶段(E-step) 中,其中的每个元素 *p*;通过以下方式计算:

$$p_{ij} = \frac{\exp\left(-\frac{1}{2\sigma^{2}} \|\boldsymbol{r}_{j}^{1} - (\boldsymbol{r}_{i}^{2} + \boldsymbol{G}_{i,*})\boldsymbol{W}\|^{2}\right)}{\left(\sum_{k=1}^{n_{2}} \exp\left(-\frac{1}{2\sigma^{2}} \|\boldsymbol{r}_{j}^{1} - (\boldsymbol{r}_{k}^{2} + \boldsymbol{G}_{k,*})\boldsymbol{W}\|^{2}\right) + \frac{w}{1-w} \frac{(2\pi\sigma^{2})^{d/2}n_{2}}{n_{1}}\right)}$$

其中W为系数矩阵,初始化为全0矩阵, $\sigma^2$ 为各向同性 斜方差矩阵,初始化值为 $\sigma^2 = -\frac{1}{dn_1n_2}\sum_{i,j=1}^{n_1,n_2} \|\boldsymbol{r}_j^1 - \boldsymbol{r}_i^2\|^2, \boldsymbol{G}$  为核对角矩阵,其元素值为 $g_{ij} = \exp(-\||\mathbf{r}_{i}^{1} - \mathbf{r}_{i}^{2}|\|^{2}/2\beta^{2})$ ,参数 $0 \le w \le 1$ ,向量 $\mathbf{r}_{i}^{1} = \mathbf{r}_{i}^{2}$ 为学习得到的节点 表示向量,分别为 $\mathbf{R}^{1'}$ 和 $\mathbf{R}^{2'}$ 的第j行和第i行,d为 向量维度, $n_{1} = n_{2}$ 为典型子图中节点的数量, $\beta$ 为大 于0的参数.最大化阶段(M-step)主要更新系数矩 阵,其计算方式如下:

 $W = (G + \lambda \sigma^2 \operatorname{diag}(P1)^{-1})^{-1} (\operatorname{diag}(P1)^{-1} PM_1 - M_2),$ 其中 diag()为对角化操作, *P* 为匹配概率矩阵, 1 为 全 1 的列向量, λ>0 为参数. 经过迭代收敛可以得 到匹配概率矩阵 *P*. 即通过 CPD 可以得到节点 *u*<sub>i</sub>与 *v*<sub>j</sub>之间的匹配概率 *p*(*u*<sub>i</sub>, *v*<sub>j</sub>) ∈ (0,1), 其中 *u*<sub>i</sub> ∈ *V*<sub>1</sub>', *v*<sub>j</sub> ∈ *V*<sub>2</sub>'.

在理想情况下,如果节点 u<sub>i</sub>与 v<sub>j</sub>存在实际的匹 配关系,那么它们之间概率应该输出为1,否则应该 输出0,概率值越大意味着两个节点的对应关系的 确定性就越大.

## 3.3 种子节点生成

本文方法为基于拓扑结构信息的无监督图对齐 算法,目前大多数无监督算法或多或少都利用到拓 扑一致性约束,对两个图分布差异较大情况下的对 齐任务均是基于有监督的,例如 Wang 等人<sup>[26]</sup>利用 双重约束机制用来平衡两个图中不同节点度差异十 分大的情况.目前大多数无监督算法都假设对齐的 两个图之间具有拓扑一致性,当两个待匹配图的拓 扑结构差异较大时,算法对齐的准确率往往较低,且 节点之间的相似性也不会太高.

为了获得高可靠的种子节点对,本文根据种子 节点的一致性概率矩阵,利用二分图最大权匹配 算法使种子节点之间的匹配概率最大化.二分图匹 配的目标是在二分图中找出一个子图,要求任意属 于该子图的两条边都不会交汇于同一个顶点,且该 子图中所有边的权值之和最大.本文使用 Kuhn-Munkres<sup>[27]</sup>算法得到种子节点集合M.本文选取两 个图中度数最大的 K 个顶点来构建典型子图是因 为其最大程度地保留了原始图的拓扑结构信息,从 而确保所得的匹配种子节点具有较高的可靠性.值 得注意的是,并非度数较大的节点一定在另一个图 中存在对应的匹配节点,为了确保种子节点的可靠 性,仅保留匹配概率大于一定阈值的种子节点对,即  $p(u_i, v_i) > \theta, (u_i, v_i) \in \mathcal{M}_s, \theta$  为可调参数,为了保障 所得种子节点的可靠性,本文取 $\theta=0.9$ ,确保所得 种子节点的准确性,以免对后续匹配造成较大的负 面影响.当两个图的拓扑结构差异较大时,匹配所得 的大部分种子节点的匹配概率会低于该阈值,若所 得种子节点的数量少于 0.5K 时,则认为源图和目标图之间存在较大的拓扑结构差异,不会进行后续的图融合匹配对齐,从而确保最终所得结果可靠有效.

## 4 基于图融合的大规模图对齐算法

在得到种子节点集*M*后,可以利用已有的有监 督图对齐算法对*G*<sub>1</sub>和*G*<sub>2</sub>进行匹配对齐.已有方法 大多先采用图表示学习算法分别获取*G*<sub>1</sub>和*G*<sub>2</sub>节点 的向量表示,然后根据种子节点将*G*<sub>1</sub>和*G*<sub>2</sub>的向量 表示转换到同一向量空间中.这一过程由于种子节 点分布与全图节点分布的不一致性可能导致转换后 的向量表示存在较大误差,影响匹配结果的准确性. 且由于*G*<sub>1</sub>和*G*<sub>2</sub>的规模较力,目前大多数的图表示 学习算法如 DeepWalk<sup>[24]</sup>、GCN<sup>[28]</sup>等均无法实现对 结构信息的快速高效提取.针对以上问题,本文提出 一种基于图融合的大规模图对齐算法,首先根据种 子节点对待匹配图融合为一个大图,然后提出一种 高效的融合图表示学习算法获得节点的向了表示, 最后利用节点向量计算它们之间的相似性并进行匹 配对齐.

#### 4.1 基于种子节点的图融合

为使已配对的高可靠种子节点在融合图中具有 更多的结构信息权重从而降低整体结构差异带来的 负面影响,本文基于已匹配的种子节点将 G<sub>1</sub>和 G<sub>2</sub> 进行融合,构造成新的大图 G<sub>3</sub>.具体地,将 G<sub>1</sub>和 G<sub>2</sub> 中已匹配的种子节点进行聚合形成一个新的节点, 其余在 G<sub>1</sub>和 G<sub>2</sub>中的节点按照原有节点之间的连接 关系添加到 G<sub>3</sub>中.

此过程可以看作将网络 G<sub>1</sub>和 G<sub>2</sub>基于种子节点 进行连通,种子节点为连接纽带.

## 4.2 融合图表示学习

目前大多数表示学习方法难以捕获图中的全局 结构信息,该问题在大规模图上尤为突出.为此,提 出一种基于个性化随机游走的图表示学习算法并应 用于 G<sub>3</sub>,可以快速得到节点的高效向量表示.个性 化随机游走能捕获到全图的结构信息,并能够高效 地表达节点之间的亲近关系,已经被广泛应用于图 数据挖掘中.具体地,对于给定 G<sub>3</sub>中的一个源节点 *u*,从节点 *u* 开始进行随机游走,每到一个节点都以 概率 *a* 停止游走并且回到源节点 *u* 重新开始,或者 以概率(1-*a*)随机选择一个邻居节点继续游走.经 过多轮游走之后,图中每个顶点被源节点 *u* 访问到 的概率会逐步收敛趋于稳定,稳定后的概率可以作为源节点 *u* 与其他节点之间的亲和度得分.把图中所有节点分别当作源节点进行个性化随机游走,可获得如下的亲和度矩阵:

$$\boldsymbol{M} = \sum_{i=0}^{\infty} \alpha (1-\alpha)^i \cdot \boldsymbol{P}^i \tag{1}$$

其中  $P = D^{-1}A$  为概率转移矩阵, A 为图的邻接矩阵, D 为对角矩阵, 对角线上元素为该节点在图中的度大小, 即 D[i,i] = [deg(i)]. 在大规模图中, 由于节点数量庞大, 直接计算个性化随机游走矩阵  $M \in \mathbb{R}^{n \times n}$ 的效率低下. 此外, 为了得到节点的低维度向量表示, 还需要对个性化随机游走矩阵进行分解以达到降维的目的. 但是个性化随机游走矩阵 M 包含着全图的高阶结构信息, M 为稠密矩阵, 对该矩阵进行奇异值分解的时间复杂度为  $O(n^3)$ , 在大规模图上会存在严重的计算效率问题, 甚至不可行.

为了提高图表示学习的可扩展性,本文提出一种近似快速分解方法,在保证有效性的前提下提高 表示学习的效率.通过分析发现,图表示学习要求输 出相似度矩阵分解后的节点低维表示,因此个性化 随机游走矩阵 M 仅为中间产物,并不要求直接计算 得出.基于这一思想本文采用 ApproxPPR<sup>[29]</sup>算法 的思想,将矩阵 M 的计算和分解集成到每一次迭代 运算中,相似度矩阵 M 可以看成是图中节点从低阶 到高除亲和度的加权和.首先对图 G<sub>3</sub>的邻接矩阵 A 进行奇异值分解:

$$A = U\Sigma V^{\mathrm{T}}$$
(2)

式中U为左奇异矩阵, $\Sigma$ 为奇异值对角矩阵,V为右 奇异矩阵.此时,可以矩阵A看作矩阵 $X_1$ 和 $Y_1$ 的乘 积,将 $X_1$ 看作节点的初始表示向量为

$$\boldsymbol{X}_1 = \boldsymbol{D}^{-1} \boldsymbol{U} \sqrt{\boldsymbol{\Sigma}} \tag{3}$$

而 $Y_1 = V \sqrt{\Sigma}$ 则为节点的上下文向量表示.初始节点 表示 $X_1$ 只包含 $G_3$ 中的一阶邻居信息,难以描述整 个网络节点间的拓扑结构.为了逐步聚合图中的高 阶结构信息,对初始嵌入向量 $X_1$ 进行迭代计算可以 得到:

$$\boldsymbol{X} = \sum_{i=1}^{\infty} \alpha (1-\alpha)^{i} \cdot \boldsymbol{P}^{i-1} \boldsymbol{X}_{1}$$
(4)

迭代后的向量表示 X 包含了图中更高阶的结构信息,其中参数  $\alpha$  的取值决定了表示向量包含的图结构信息的范围,  $\alpha$  取值过低则表示向量能够包含更多的高阶信息,但需要较多的迭代次数才能收敛;取值过低则会包含更多的局部结构信息,比如在极端情况下  $\alpha$  取值为 1 时,则学习得到的表示向量只包

含节点的一阶邻居结构信息.为了综合算法的有效 性和时间效率,本文中α取值为 0.15.

由于邻接矩阵A为稀疏矩阵,可以采用BKSVD<sup>[30]</sup> 进行近似奇异值分解得到A≈UΣV<sup>T</sup>的近似解.该奇 异值分解的时间复杂度与网络中边的数量成线性关 系,远低于分解稠密矩阵 M 需要的O(n<sup>3</sup>)复杂度, 大大提高了图表示学习方法的可扩展性.

最终得到的图表示向量为 *R*=[*X*⊕*Y*],⊕为对 表示向量 *X* 和 *Y* 进行拼接操作.

4.3 基于社区结构的高阶信息快速传播

在对式(4)进行高阶信息计算与聚合过程中,需 要大量的大规模矩阵乘法迭代运算,导致计算效率 依旧不高.为了提高这一部分的计算效率,本文利用 图中普遍存在的社区结构对高阶信息传播过程进行 加速计算.在真实世界形成的图上往往会存在一定 的聚集现象[31],也称为社区结构,在同一社区内的 节点之间连接关系更加紧密,而属于不同社区的节 点其连接关系相对而言比较稀疏[32-35],社区结构在 真实图数据中广泛存在,如社交网络中具有其同兴 趣爱好的用户之间会有更多的连接关系,从而形成 兴趣社团[36];而在生物医药领域,研究人员往往将 具有相似生物功能的蛋白质和药物用边连接起来, 从而会形成生物蛋白质网络中的社区结构[37].如图 3 所示. 图中出现两个社区结构, 可以明显看出社区 内的节点相比社区间的节点之间的连接更加紧密. 信息的传播依赖于节点之间的连接,其连接密集程 度对传播效率会产生极大的影响,因此节点信息在 社区内与社区间的传播存在较大差异,本文利用图 中存在社区结构这一特点对高阶信息传播进行加速 计算.



图 3 社区聚集现象

在式(4)中为了得到图的表示向量 X,需要大量 迭代使图中的高阶信息融入初始表示向量 X<sub>1</sub>中,该 过程可以看作信息在图上不断的传播和聚合.由于 图中存在社区结构,社区内的节点之间连接较为紧 密,导致节点信息在其相应的社区内反复传播,难以 扩散到社区外,只有经过一定次数的迭代传播后,节 点信息才能逐步扩散到全图上<sup>[38]</sup>.因此,该信息在 图中的传播过程可以分为社区传播和逃逸传播 2 个 阶段,基于发现可以对表示向量 X 进行快速计算. 4.3.1 社区传播阶段

假定社区传播阶段的迭代次数为 s,在对社区 传播阶段进行计算时,首先对式(4)中的前 *l*<s 项 进行求和计算,可以得到 **X**<sub>Reality</sub>:

$$\boldsymbol{X}_{\text{Reality}} = \sum_{i=1}^{l} \alpha (1-\alpha)^{i} \cdot \boldsymbol{P}^{i-1} \boldsymbol{X}_{1}$$
(5)

X<sub>Reality</sub>包括了初始表示向量 X<sub>1</sub>前 *l* 次迭代,其计算 结果为社区传播阶段中前 *l* 次节点信息在网络中传 播得到的向量表示,其中 *l* 为可调参数.设 X<sub>Community</sub> 为整个社区传播阶段得到的向量矩阵:

$$\boldsymbol{X}_{\text{Community}} = \sum_{i=1}^{s} \alpha (1-\alpha)^{i} \cdot \boldsymbol{P}^{i-1} \boldsymbol{X}_{1}$$
 (6)

由于图中存在社区结构,因此在一定次数内的信息 传播其实是在社区内反复进行的,属于社区内的节 点大概率会将其信息传播到同一社区内的节点上,在 前 *l* 次传播中获得信息的节点将在接下来的(*s*-*l*)次 迭代中有很大概率会再次获得该信息,因此在社区 传播阶段接下来的第 *l*+1 次迭代到第 *s* 次迭代,节 点获得的信息可以看作与前 *l* 次迭代所得信息成比 例,即 *X*<sub>Community</sub>可以从 *X*<sub>Reality</sub>中计算得到:

$$X_{\text{Community}} \approx \delta \cdot X_{\text{Reality}}$$
(7)

式中 δ 为比例系数,对比前 *l* 次迭代与整个社区传播阶段的总得分可以求得比例系数 δ.

$$\delta = \frac{\left\|\sum_{i=1}^{s} \alpha (1-\alpha)^{i} \cdot \boldsymbol{P}^{i-1}\right\|_{1}}{\left\|\sum_{i=1}^{l} \alpha (1-\alpha)^{i} \cdot \boldsymbol{P}^{i-1}\right\|_{1}}$$
(8)

其中 *l*<*s*,由于概率转移矩阵 *P* 的每一项都是非负 实数,可以将式(8)如下表示:

$$\delta = \frac{\sum_{i=1}^{3} \alpha (1-\alpha)^{i} \| \boldsymbol{P}^{i-1} \|_{1}}{\sum_{i=1}^{l} \alpha (1-\alpha)^{i} \| \boldsymbol{P}^{i-1} \|_{1}}$$
(9)

概率转移矩阵 P 为行随机矩阵,其每一行元素之和 为 1,可以得出  $P^{i-1}$ 也为行随机矩阵,且  $\|P^{i-1}\|_1 = n$ ,由此可以得出  $\delta$  的表达式为

$$\delta = \frac{\sum_{i=1}^{s} \alpha (1-\alpha)^{i}}{\sum_{i=1}^{l} \alpha (1-\alpha)^{i}} = \frac{1-(1-\alpha)^{s}}{1-(1-\alpha)^{l}}$$
(10)

基于以上计算可以快速得到初始表示向量 X<sub>1</sub> 在社 区传播阶段进行高阶信息聚合后的空间向量矩阵  $X_{\text{Community}}$ .

4.3.2 逃逸传播阶段

在得到 *X*<sub>Community</sub>后,信息的传播将不再受限制 于图中的社区结构,而将进一步扩散到全图的节点 上,此阶段称为逃逸阶段,设经过逃逸阶段传播后的 表示向量矩阵为 *X*<sub>Eescape</sub>.在社区传播阶段,距离源节 点越近,节点从源节点获得的信息将会越多.而在逃 逸阶段由于经过多次迭代,在系数 0<1-α<0 的影 响下,此时图中节点经过高阶传播后信息与源节点 的相关性不断被削减,即在逃逸阶段节点所获得 的信息与到传播源的远近无关,此时,度越大的节 点越有机会接收到邻居节点传播的信息,并且所获 得的信息与邻居节点所拥有信息正相关,此思想与 PageRank 的基本假设一致,

因此可以用原始的 PageRank 算法对逃逸阶段 的信息传播进行计算.首先计算图G<sub>3</sub>的 PageRank 值, 得到每个节点的权重,算法的累计幂达代形式如下:

$$pagerank(G) = \sum_{i=0}^{\infty} \alpha (1-\alpha)^{i} e \mathbf{P}^{i}$$
(11)

其中  $e \in R^{1 \times n}$ ,其中每一项元素值都为 $\frac{1}{n}$ ,加入网络 中节点的个数,最终聚合后的表示向量为

$$\boldsymbol{x}_{\text{Eescape}}^{i} = \sum_{i=s+1}^{\infty} \alpha (1-\alpha)^{i} \boldsymbol{e} \boldsymbol{P}^{i} \boldsymbol{X}_{1}$$
(12)

当 $\|\mathbf{x}_{\text{Eescape}}^{i+1} - \mathbf{x}_{\text{Eescape}}^{i}\|_{1} \leq 10^{-6}$ 时停止迭代,最终得 到表示矩阵  $\mathbf{X} \approx \mathbf{X}_{\text{Community}} + n e^{T} \mathbf{x}_{\text{Eescape}}$ .值得注意的 是,式(12)中 e 为行向量,因此在计算时向量 e 先与 矩阵 P 相乘得到新的行向量再继续迭代运算,而不 是先计算矩阵 P 的多次幂,从而避免了大型矩阵的 乘法运算.

#### 4.4 基于节点相似度的快速匹配对齐

在得到节点的向量表示后,可以计算图  $G_1$ 和  $G_2$ 中所有节点两两之间的相似性,然后按照 Kuhn-Munkres<sup>[27]</sup>算法进行匹配.但由于节点数量庞大,该 算法计算效率较低,为此,本文采用一种快速的匹配 对齐算法进行计算.首先,计算节点表示向量在空间 中的欧氏距离来描述它们之间的相似性,假设节点  $u \in G_1, v \in G_2, 那么它们之间距离为d_{u,v} = \|v_u - v_v\|_2,$ 对于每一个属于  $G_2$ 的节点,在  $G_1$ 中寻找与其距离 最近的节点进行匹配,一旦匹配则移除该节点,该算 法虽牺牲了一定的准确性,但相比 Kuhn-Munkres 算法而言,其时间复杂度从  $O(n^3)$ 降为了  $O(n^2),$ 提 高了计算效率.

值得注意的是,本文核心算法虽然间接借鉴了

ApproxPPR 算法的思想和直接采用 Kuhn-MUnkres 等算法,但本文的核心贡献在于提出一种基于种子 节点融合的大规模图对齐算法,与以往基于图表示 学习的图对齐方法不同,本文所提方法直接在融合 图上进行无监督的图表示学习,避免了表示空间的 转换,从而大大提高了匹配准确率.此外,为了提高 图表示学习算法在大规模融合图上的计算效率,本 文算法利用图的社区结构来加速高阶结构信息的传 播,在保证有效性的同时大大提高了图表示学习的 计算效率,使得本文算法能够较好地解决大规模图 对齐问题.

## 5 实验与结果分析

#### 5.1 实验数据与运行环境

本文在6个被广泛使用的公开数据集上进行了 实验,表1列出了数据集的详细信息. PPI<sup>[39]</sup>为人类 蛋白质网络,BlogCatalog<sup>[39]</sup>为社交网络,由不同博 主与其社会关系构成,Enron<sup>[40]</sup>是由联邦能源监管 委员会在调查期间公布并发布的邮件通信网络, Slashdot<sup>[41]</sup>为资讯科技网络,其所有新闻由用户提供, 用户之间可以形成朋友或者敌人关系,Youtube<sup>[42]</sup>为 视频共享网站,用户之间可以建立友谊关系并进行 视频分享, Pokec<sup>[43]</sup>是斯洛伐克最受欢迎的在线社 交网络 所有的数据集都可以在公开网站或者论文 中获得.本文主要解决基于拓扑结构信息的无监督 图对齐问题,由于算法仅依赖图的拓扑结构,参照以 往类似算法的实验假设[44-46]:待匹配的两个图往往 共享共同的连接关系 例如微信社交网络与 QQ 好 友关系网络都是以用户之间真实的熟人关系为基础 的,但由于人类活动的随机性和数据采集不完全等 原因会导致此2个社交网络的拓扑结构具有一定的 差异性,但都可以看作是从用户真实的社会网络按 照一定的比例采样连接关系得到的,因此,在实验 时,将原始图表示为G,对原始图G随机删除一定比 例的边后得到最大连通子图  $G_1$ 和  $G_2$ ,设删除边的 比例为 $\gamma$ ,然后对 $G_1$ 和 $G_2$ 进行对齐.此实验设置符 合大多数的图对齐任务的应用场景.由于随机裁剪 边生成待匹配图具有一定的不确定性,为了获得更 加稳定的实验结果,所有实验重复执行 20 次.从运 行结果来看,随机裁剪对运行结果并没有产生明显 的影响,因此该操作具备普适性,在所有实验中仅使 用图结构信息来完成对齐任务,且初始时均不包含 任何先验的匹配信息.

表 1 实验数据集				
数据集名称	类型	节点数量*	边数量	
PPI	蛋白质网络	3 8 9 0	76584	
BlogCatalg	社交网络	10312	333983	
Enron	邮件通信网络	36692	183831	
Slashdot	新闻网络	77360	905468	
Youtube	视频共享网络	1134890	2987624	
Pokec	社交网络	1632803	30622564	

为了验证本文所提图对齐算法的有效性,与7种 目前被广泛使用的无监督图对齐算法进行对比,分别 是BigAlign<sup>[12]</sup>、IsoRank<sup>[13]</sup>、FINAL<sup>[14]</sup>、REGAL<sup>[15]</sup>、 GAlign<sup>[17]</sup>、WAlign<sup>[18]</sup>和NAWAL<sup>[22]</sup>.本文所有实 验代码均运行在系统为Unbuntu 18.04、CPU为 Intel(R) Xeon(R) Gold 6266C CPU@3.00GHz和 内存为128GB的服务器上.表2列出了本文算法的 参数设置表.在分析不同参数对模型效果的影响时, 为控制实验变量,仅改变待分析参数的取值,其余参 数按照表2中的默认值进行设置并且保持不变.本 文以图对齐匹配准确率作为不同算法的评价指标, 在计算时将匹配概率最大的节点对作为对充节点, 计算正确对齐节点占图*G*1中总节点数量的5.7比. 所有结果均为20次实验的平均值.

表 2 实验参数设置表

参数名称	参数描述	取值范围	默认值
α	个性化随机游走的停止概率	$0 \sim 1.0$	0.15
\$	社区传播迭代次数	$6 \sim 18$	9
l	社区传播初始迭代次数	$2 \sim 7$	4
β	典型子图节点数量调节系数	$2 \sim 12$	4
d	嵌入向量的维度	$40 \sim \! 18$	80
θ	种子节点匹配概率阈值	$0 \sim 1.0$	0.90

#### 5.2 实验结果及分析

5.2.1 种子节点匹配结果分析

图 4 为本文所提算法在种子节点匹配时的实验 结果.实验结果表明本文方法得到的种子节点具有 较高的匹配准确度和可靠性,虽然随着移除边比例 的增加,种子节点的匹配准确度也在逐步下降.其中



Pokec 数据集上准确性下降最快,但即使在移除边的比例达到 20%的情况下依然取得了 0.92 以上的匹配准确率,确保了种子节点的可靠性,为后续基于图融合的匹配奠定了基础.

5.2.2 全图匹配结果分析

图 5 为在 6 个数据集上各算法的实验结果对比. 本文所提方法在 PPI 和 BlogCatalog 两个数据集上 在所有移除边比例相同的情况下均取得了最好的匹 配准确率.与其他方法相比,本文所提方法的性能随 移除边比例增加呈缓慢下降趋势,这是因为随着移 除边比例的增加,待匹配图中包含的结构信息会随 之减少,匹配准确率有所下降也在情理之中.与本文 方法形成鲜明对比的是,其他对比方法的匹配准确 率随移除边比例的增加呈现急剧下降的趋势.特别 是 REGAL 与 WAlign 2 个方法,当移除边的比例 不高时,待匹配图的拓扑结构较为相似,二者的匹配 准确率除低于本文所提算法外,远高于其他 3 种方 法.但当移除边的比例大于 15%时,二者的匹配准 确率急剧下降,其匹配准确率几乎与其他 3 种方法 持平.

值得注意的是,当移除边的比例小于3%时, REGAL 方法在数据集 Enron、Slashdot 和 Youtube 上的匹配准确率要高于本文所提算法, WAlign 方 法在数据集 Enron 上的匹配准确率高于本文所提 出的算法.得到该匹配结果是因为 REGAL 方法通 过计算待匹配节点邻居的度大小信息,在待匹配图 结构差异不大的情况下实际匹配节点具有相似的结 构信息,学习到的表示向量在特征空间上更加接近, 从而能够取得较高准确率. 而 WAlign 方法最小化 不同图节点表示向量之间的 Wasserstein 距离找到 最优对齐结果,相当于将问题转化为最小代价的二 分图最优匹配问题,在移除边比例较小的情况下能 够使节点表示向量在不同图中具有相似的概率分 布,从而达到较高的图对齐准确率,然而 WAlign 方 法随着图规模的扩大,其计算所耗费时间会急剧上 升.同时这两种方法对待匹配图的拓扑结构差异性 较为敏感,二者的性能随移除边比例的增加会急剧 降低.

图 5(a)到图 5(d)的结果表明,NAWAL 模型 较其他方法取得了相对稳定的匹配准确率,其性能 未出现随着删除边比例的增加而急剧下降的情况. 这是因为 NAWAL 采用对抗训练的方式学习源图 和目标图之间的映射关系,对两个图之间的结构差 异变化具有较高的适应性. 然而 NAWAL 模型的准



确率远低于本文所提方法.此外,GAlign 模型在PPI、 BlogCatalog、Enron 上取得了很低的匹配准确率, 这是因为该模型在节点对齐时非常依赖节点的属性 信息,而本实验中仅使用节点的结构信息,导致了较 高的错误匹配率.

特别说明,在数据集 Slashdot 上 WAlign、GAlign 等方法运行时间大于三天而无法得出匹配结果,在数 据集 Youtube 上,除了本文提出的方法和 REGAL 方法,其他对比方法皆因效率或者内存原因而无法 完成对齐实验,因此得到如图 5(d)和(e)所示的实 验结果.在拥有 3 千多万条边的大型数据集 Pokec 上,只有本文提出的算法在三天内得到匹配结果且 不超出内存限制,实验结果如图 5(f)所示.

#### 5.3 参数对实验影响

本小节验证信息传播参数 s 和 l、种子节点参数 β、向量表示维度 d 等参数对实验结果的影响.其中 参数 s 和 l 分别控制着社区传播阶段和逃逸传播阶 段的迭代次数,二者的大小会对表示向量 X 计算的 准确度与计算效率造成影响,将初始化向量 X 进行 高阶信息传播所造成的计算误差定义为

 $error = \| \mathbf{X}_{\text{Community}} + \mathbf{X}_{\text{Eescape}} - \sum_{i=1}^{\infty} \alpha (1-\alpha)^{i} \cdot \mathbf{P}^{i-1} \mathbf{X}_{1} \|_{1}.$ 在研究参数 *s* 对 *error* 的影响时,为了控制变量,将 *l* 值

在研究参数 s 約 error 的影响时, 为 ] 控制变重,将 l 值 固定为4,改变参数 s 的大小,在大型数据集 Youtube 上进行实验,结果如图 6(a)所示.实验结果表明计 算误差 error 在初始阶段会随着参数 s 的增大而减 少,但当s值增长到9左右时,误差开始随着参数s 增大.基于社区结构进行高阶信息传播所产生的误 差由社区传播  $X_{\text{Community}}$  和逃逸传播  $X_{\text{Escape}}$  构成,随 着参数》的增大,社区传播阶段高阶信息的聚合需 要估计的迭代次数会增多,从而造成此阶段产生 的误差增长, 而逃逸传播阶段需要估计的迭代次 数会减少\因此误差减小,在 s 刚开始增长时逃逸 阶段的误差减小量大于社区传播阶段误差增长 量,而随着s的进一步增大,误差的增长量会逐渐大 于减小量,从而造成整体误差随着参数。的增长先 减小继而再上升的现象.根据实验结果,本文设置参 数 s=9. 图 6 表明参数 s 和 l 都会对误差造成影响, 其中参数 l 会对参数 s 的拐点值的产生造成影响, 参数 s 对 l 无影响. 此外,由式(10)可知,比例系数 δ 可由 s 的值直接计算得到,此部分计算时间与 s 的取 值大小无关,因此本文所提算法的运行时间主要由 参数 l 的取值大小决定,因此在图 6(a)中未画出运 行时间.

在研究参数 l 对误差与计算效率的影响时,将 s 值固定设置为 9,改变 l 值大小,实验结果如图 6(b) 所示.随着 l 值的增长,社区传播阶段对向量 X<sub>1</sub>的 传播迭代次数增加,因此耗费的时间会随之增长,但 会减小该阶段带来的计算误差,考虑到准确性和运 行时间之间的平衡,本文设置 l=4.





本文通过改变典型子图的大小来验证预匹配过 程中获得的种子数量对算法性能的影响,由于典型 子图的节点数量  $K = \beta \sqrt[2]{n}$ ,其中 n 为原始图节点的数 量, $\beta$  为可调参数,图 7 展示了 3 的不同取值在 PPI 和 BlogCatalg 两个数据集上的实验结果.当  $\beta \leq 4$  时, 图对齐的准确率会随着  $\beta$  的增大急剧提升,但当  $\beta > 4$  时,准确率的提升明显放缓,而种子节点匹配 所耗费的时间会随着  $\beta$  的增大而线性上升.以上结 果表明,当典型子图过小时,会导致预匹配场得的种 子节点数量较少,从而导致后续基于融合图的匹配 结果准确度较低.而当预匹配获取到一定数量的种 子节点后,增加更多的种子节点也不会明显提高基 于融合图的匹配准确率,且典型子图过大会导致预



匹配阶段消耗的时间大大增加,为了兼顾准确性和 计算效率,本文设置 $\beta=8$ .

为了选择合适的维度值,本节将图表示学习方 法输出的节点向量维度设置从 40 增加到 180,在 BlogCatalg 和 Enron 两个数据集上进行实验,评估 向量维度对图对齐方法准确率和整体效率的影响, 实验结果如图 8 所示.准确率随着维度 d 的增加会 缓慢提升,但所消耗的计算时间也会随之增加,且增 长速度远高于准确率的提升.这主要是由于本文所 提的图对齐方法利用表示向量之间的相似性实现节 点对齐这一步骤会随着维度增加而需要更长的搜索 时间开销.综合考虑算法准确率和计算效率,本文设 置表太向量的维度 d=80.









图 8 表示向量维度 d 对准确率和效率影响

7 期

#### 5.4 计算效率对比

表 3 详细描述了不同方法在实验数据集上的运 行时间,本文所提的图对齐方法在计算效率上明显 优于其他方法,与传统图对齐方法 BigAlign、FINAL 和 IsoRank 相对比,本文所提方法 LGA 的运行时 间在所有数据集上均处于绝对优势.与同样使用表 示学习的图对齐方法 REGAL 相比,REGAL 以在 结构差异较大的图对齐任务中具有一定的鲁棒性为 代价,避免了基于随机游走采样的计算开销,因此效 率取得一定的提升,但与本文的方法 LGA 相对比 依然有较高的计算成本,特别是在大规模图数据集 上,其计算时间开销会显著增加,REGAL 方法在 Youtube 上运行时间将近 30 个小时,在拥有 3 千万 条边的大型数据集 Pokec 上无法在三天之内完成计 算.与使用图卷积网络的 Walign 方法相比,WAlign

>!

在小规模数据集上有较高的时间效率优势,但随着 节点数量的增长,其计算损失函数与参数更新会耗 费大量的时间和内存,在接近 8 万个节点数据集上 Slashdot上运行时间大于三天.图 9 为各方法运行 时间与节点数量之间的关系曲线图,随着节点数量 的增长,FINAL、IsoRank 和 BigAlign 的运行时间 会急剧增加.而 REGAL、WAlign、LGA 三种方法在 节点数量大于 35k 后时间消耗也开始大幅度增加, 在大规模数据集上不具备可扩展性.与这些方法不 同的时,本文方法 LGA 随着网络节点数量的增加 其运行时间始终保持平缓增长,在不同规模的数据 集上均能高效地实现图对齐.特别是本文方法在待 匹配图之间的结构差异较大时仍能够达到良好的节 点匹配准确率,并且在运行效率方面相比其他方法 有较大的优势.

$\mathbb{Z}$ $\mathbb{X}$ $\mathbb{R}$ $3$ 不同方法的运行时间比较 (单位: s					(単位:s)	
名称	BigAlign	FINAL	IsoRank	REGAL	WAlign	LGA
PPI	43	108	62	14	3	13
BlogCatalg	452	678	543	218	24	34
Enron	2823	16155	9537	165	1223	73
Slahdot	21 4 4 3	59370	34 301	1068	时间>3天	104
Youtube		内存溢出	1	107 602	时间>3天	11031
Pokec		内存溢出		时间>3天	内存溢出	27 036
60 000 50 000 40 000 30 000 20 000	LGA REGAL IsoRank FINAL WAlign		机 構 ) 于 オ	在进行图对齐住 • 相同的效果, · 进行信息传播.	壬务时能达到与直 并且得到的匹配 ⅠⅠG/	重接迭代信息传 是准确率明显高 A-Comm



#### 5.5 消融实验

本节利用消融实验来验证本文提出的基于社区 结构进行高阶信息传播(LGA-Comm)这一机制在 图表示学习时的有效性,在表示学习阶段设置直接 迭代进行高阶信息传播(LGA-Full)和不进行信息 传播(LGA-Adj),在 PPI、BlogCatalog、Enron 和 Youtube 四个数据集上进行图对齐任务并进行对 比,实验过程中将移除边的比例设置为 10%,基于 不同传播方式取得的匹配准确率如图 10 所示.从图 中可以看出,基于社区结构进行高阶信息传播这一



本节对基于社区结构进行高阶信息传播和直接 迭代信息传播的效率也进行了比较,结果如表 4 所 示,基于社区结构进行信息传播耗时不到直接迭代 计算的十分之一,具有较大的效率优势.相比不进行 信息传播,节点表示向量只能学习到图网络的局部 结构信息,使得节点向量在空间中推理能力下降,其 最终的准确率低于另外两种方法.

名称	LGA-Comm	LGA-Adj
PPI	0.87	0.07
BlogCatalog	4.92	0.24
Enron	5.18	0.41
Youtube	59.48	4.28

因此,本实验表明图数据的全局结构信息能明 显提高图对齐任务的准确率,并且本文提出的基于 社区结构进行高阶信息传播在保证有效性的情况下 能够大幅度降低学习图全局结构信息所需耗费的 时间.

## 6 结 论

为了解决大规模图对齐任务中的初始种子节点 发现难和图表示学习算法效率与有效性难以兼顾的 问题,本文提出了一种基于稀疏矩阵分解和个性化 随机游走相结合的图表示学习方法,能够对大规模 网络进行高效学习,得到的表示向量能够捕获图的 全局结构信息,很好地适用于图对齐任务,在此基础 上提出一种基于种子节点融合的大规模图对齐算法 LGA,与以往利用种子节点监督图表示学习的范式 不同,本文所提方法直接在融合图上进行无监督的 图表示学习,直接将待匹配的图数据映射到统一的 向量空间中,避免了由于表示空间转换导致的匹配 准确率较低的问题.最后在多种类型和不同规模的 数据上进行了广泛的实验验证,结果表明本文所提 方法的能够取得远高于其他方法的匹配准确率和计 算效率,能够应用于大规模图数据的高效对齐.

本文方法完全利用图的拓扑结构信息来实现图 对齐,具有较广的应用范围.然而,部分图数据包含 了丰富的属性信息,对于提高图对齐的准确性具有 重要作用,为此,在今后工作中将继续改进本文算 法以便结合网络的属性信息进一步提高图对齐的准 确率.

#### 参考文献

- Marin A, Wellman B. Social network analysis: An introduction. The SAGE Handbook of Social Network Analysis, SAGE, 2011; 11-25
- [2] Perc M. Growth and structure of Slovenia's scientific collaboration network. Journal of Informetrics, 2010, 4(4): 475-482
- [3] Alm E, Arkin A P. Biological networks. Current Opinion in Structural Biology, 2003, 13(2): 193-202

- [4] Liu X, Murata T, Kim K S, et al. A general view for network embedding as matrix factorization//Proceedings of the 12th ACM International Conference on Web Search and Data Mining (WSDM'19). New York, USA, 2019; 375-383
- [5] Yin Y, Wei Z. Scalable graph embeddings via sparse transpose proximities//Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining (KDD'19). New York, USA, 2019: 1429-1437
- [6] Man T, Shen H, Liu S, et al. Predict anchor links across social networks via an embedding approach//Proceedings of the 25th International Joint Conference on Artificial Intelligence (IJCAI'16). Palo Alto, USA, 2016: 1823-1829
- [7] Yu Liang, Gao Lin, Sun Peng-Gang. Research on algorithms for complexes and functional modules prediction in proteinprotein interaction networks. Chinese Journal of Computers, 2011, 34(7): 1239-1251(in Chinese)

(鱼亮,高琳,孙鹏岗.蛋白质网络中复合体和功能模块预测 算法研究.计算机学报,2011,34(7):1239-1251)

[8] Wu An-Biao, Yuan Ye, Qiao Bai-You, et al. The influence maximization problem based on large-scale temporal graph. Chinese Journal of Computers, 2019, 42(12): 2647-2664(in Chinese)

(吴安彪,袁野,乔百友等.大规模时序图影响力最大化的算法研究.计算机学报,2019,42(12):2647-2664)

 9] Zhang Fu, Yang Lin-Yan, Li Jian-Wei, et al. An overview of entity alignment methods. Chinese Journal of Computers, 2022, 45(6): 1195-1225(in Chinese)

**张**富,杨琳艳,李健伟等.实体对齐研究综述.计算机学 报,2022,45(6):1195-1225)

- [10] Liu L, Chome W K, Li X, et al. Aligning users across social networks using network embedding//Proceedings of the 25th International Joint Conference on Artificial Intelligence (IJCAI'16). Palo Alto, USA, 2016; 1774-1780
- [11] Klau G W. A new graph-based method for pairwise global network alignment. BMC Bioinformatics, 2009, 10(1): 1-9
- [12] Koutra D, Tong H, Lubensky D. Big-Align: Fast bipartite graph alignment//Proceedings of the 2013 IEEE 13th International Conference on Data Mining. Dallas, USA, 2013: 389-398
- [13] Singh R, Xu J, Berger B. Global alignment of multiple protein interaction networks with application to functional orthology detection. Proceedings of the National Academy of Sciences, 2008, 105(35): 12763-12768
- [14] Zhang S, Tong H. FINAL: Fast attributed network alignment //Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD'16). New York, USA, 2016: 1345-1354
- [15] Heimann M, Shen H, Safavi T, et al. REGAL: Representation learning-based graph alignment//Proceedings of the 27th ACM International Conference on Information and Knowledge Management (CIKM'18). New York, USA, 2018: 117-126

- [16] Chu X, Fan X, Yao D, et al. Cross-network embedding for multi-network alignment//Proceedings of the World Wide Web Conference (WWW'19). New York, USA, 2019: 273-284
- [17] Trung H T, Van Vinh T, Tam N T, et al. Adaptive network alignment with unsupervised and multi-order convolutional networks//Proceedings of the 2020 IEEE 36th International Conference on Data Engineering (ICDE). Dallas, USA, 2020; 85-96
- [18] Gao J, Huang X, Li J. Unsupervised graph alignment with Wasserstein distance discriminator//Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining (KDD'21). New York, USA, 2021: 426-435
- [19] Wang C, Zhao Z, Wang Y, et al. DeepMatching: A structural seed identification framework for social network alignment// Proceedings of the 2018 IEEE 38th International Conference on Distributed Computing Systems (ICDCS). Vienna, Austria, 2018: 600-610
- [20] Zhou F, Liu L, Zhang K, et al. DeepLink: A deep learning approach for user identity linkage//Proceedings of the IEEE Conference on Computer Communications (IEEE INFOCOM 2018). Honolulu, USA, 2018; 1313-1321
- [21] Derr T, Karimi H, Liu X, et al. Deep adversaria network alignment//Proceedings of the 30th ACM International Conference on Information & Knowledge Management (CIKM'21). New York, USA, 2021: 352-361
- [22] Nguyen T T, Pham M T, Nguyen T T, et al. Structural representation learning for network alignment with selfsupervised anchor links. Expert Systems with Applications, 2021, 165: 113857
- [23] Zheng C, Pan L, Wu P. JORA: Weakly supervised user identity linkage via jointly learning to represent and align. IEEE Transactions on Neural Networks and Learning Systems, 2022, (1): 1-12
- [24] Perozzi B, Al-Rfou R, Skiena S. DeepWalk: Online learning of social representations//Proceedings of the 20th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD'14). New York, USA, 2014: 701-710
- [25] Hirose O. A Bayesian formulation of coherent point drift. IEEE Transactions on Pattern Analysis and Machine Intelligence, 2020, 43(7): 2269-2286
- [26] Wang Y, Peng Q, Wang W, et al. Network alignment enhanced via modeling heterogeneity of anchor nodes. Knowledge-Based Systems, 2021, 250: 109116
- [27] Zhu H, Zhou M C, Alkins R. Group role assignment via a Kuhn-Munkres algorithm-based solution. IEEE Transactions on Systems, Man, and Cybernetics-Part A: Systems and Humans, 2011, 42(3): 739-750
- [28] Chiang W L, Liu X, Si S, et al. Cluster-GCN: An efficient algorithm for training deep and large graph convolutional networks//Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining (KDD'19). New York, USA, 2019: 257-266

- [29] Yang R, Shi J, Xiao X, et al. Homogeneous network embedding for massive graphs via reweighted personalized PageRank. Proceedings of the VLDB Endowment, 2020, 13(5): 670-683
- [30] Musco C, Musco C. Randomized block Krylov methods for stronger and faster approximate singular value decomposition. Advances in Neural Information Processing Systems, 2015, 2015; 1396-1404
- [31] Li Y, Wang Y, Zhang T, et al. Learning network embedding with community structural information//Proceedings of the 28th International Joint Conference on Artificial Intelligence. Macao, China, 2019: 2937-2943
- [32] Girvan M, Newman M E J. Community structure in social and biological networks. Proceedings of the National Academy of Sciences, 2002, 99(12): 7821-7826
- [33] Newman M E J. Detecting community structure in networks. The European Physical Journal B, 2004, 38(2): 321-330
- [34] Blondel V D, Guillaume J L, Lambiotte R, et al. Fast unfolding of communities in large networks. Journal of Statistical Mechanics: Theory and Experiment, 2008, 2008 (10): P10008
- [35] Benson A R, Gleich D F, Leskovec J. Higher-order organization of complex networks. Science, 2016, 353(6295): 163-166
- [36] Steinhaeuser K, Chawla N V. Community detection in a large real-world social network//Social Computing, Behavioral Modeling, and Prediction. Boston, USA: Springer, 2008.
   108, 175
- [37] Juhr, Singh L O, Clauset A, et al. Exploring community structure in biological networks with random graphs. BMC Bioinformatics, 2014, 15(1): 1-14
- [38] Yoon M, hung J, Kang U. TPA: Fast, scalable, and accurate method for approximate random walk with restart on billion scale graphs Proceedings of the 2018 IEEE 34th International Conference on Data Engineering (ICDE'18). Paris, France, 2018; 1132-1143
- [39] https://github.com/THUDM/ProNE/tree/master/data
- [40] http://snap. stanford. edu/data/email-Enron. html
- [41] http://snap. stanford. edu/data/soc-Slashdot0811. html
- [42] http://snap. stanford. edu/data/com-Youtube. html
  - [43] http://snap. stanford. edu/data/soc-Pokec. html
  - [44] Pedarsani P, Figueiredo D R, Grossglauser M. A Bayesian method for matching two similar graphs without seeds//Proceedings of the 2013 51st Annual Allerton Conference on Communication, Control, and Computing(Allerton). Monticello, USA, 2013: 1598-1607
  - [45] Zhang S, Tong H, Tang J, et al. Incomplete network alignment: Problem definitions and fast solutions. ACM Transactions on Knowledge Discovery from Data (TKDD), 2020, 14(4): 1-26
  - [46] Wang C, Wang Y, Zhao Z, et al. Credible seed identification for large-scale structural network alignment. Data Mining and Knowledge Discovery, 2020, 34(6): 1744-1776



WANG Chen-Xu, Ph. D., associate professor. His current research interests include graph data mining, artificial intelligence and network security. **ZHOU Jun-Ming**, M. S. His current research interests include graph alignment and network representation learning.

**JIANG Pei-Jing**, M. S. candidate. His current research interests include graph alignment and network representation learning.

#### Background

The topic studied in this paper is graph alignment in the field of graph data mining. Graph alignment has long been studied because of its important applications in various fields such as social network analysis, bioinformatics, and computer vision. At present, researchers focus on employing graph representation learning models for graph alignment with development of artificial intelligence.

In this paper, we propose a large scale unsupervised graph alignment framework based on topological structure representation learning. We first select a typical subgraph from each of the graphs as a candidate set of seed nodes. The local topological information is used to retrieve a set of highly reliable seed node pairs. Then, we use the seed nodes to fuse the matching graphs, and propose an efficient unsupervised representation learning algorithm to map the fused graph into a unified vector space. Finally, large-scale graph alignment is realized based on the learned node vectors. Compared with existing methods, the proposed approach uses the least time in large-scale graph alignment tasks and achieves the best performance of alignment accuracy. Moreover, the structural differences of graphs have limited impacts on the performance of the proposed method.

The research presented in this paper is supported in part by the National Natural Science Foundation of China (No. 62272379), the Natural Science Basic Research Plan in Shaanxi Province (No. 2021JM-018), the National Key R&D Program of China (No. 2021YFB1715600), and the Fundamental Research Funds for the Central Universities (No. 1191320006).