

有向加权网络中的改进 SDNE 算法



马扬 程光权 梁星星 李妍 杨雨灵 刘忠

国防科技大学系统工程学院 长沙 410073

(yang_ma_cn@163.com)

摘要 网络化的数据形式能够表示实体以及实体和实体之间的联系,网络结构在现实世界中普遍存在。研究网络中节点和边的关系具有重要意义。网络表示技术将网络的结构信息转换为节点向量,能够降低图表示的复杂度,同时能够有效运用到分类、网络重构和链路预测等任务中,具有很广泛的应用前景。近年提出的 SDNE(Structural Deep Network Embedding)算法在图自编码领域取得了突出成果,文中针对网络表示算法 SDNE 在有向、有向网络中的局限性,从网络结构和衡量指标两个角度入手,提出了新的基于图自编码的网络表示模型,在原有节点向量的基础上引入了接收向量和发出向量的概念,优化了自编码器的解码部分,进而优化了神经网络的结构,减少了网络的参数以加快收敛速度;提出了基于节点度的衡量指标,将网络的加权特性反映在网络表示的结果中。在 3 个有向加权数据集中的实验证明,在进行网络重构和链路预测任务时,所提方法能够取得比传统方法和 SDNE 原始方法更好的结果。

关键词: 复杂网络;网络表示;网络重构;链路预测;自编码

中图分类号 TP311

Improved SDNE in Weighted Directed Network

MA Yang, CHENG Guang-quan, LIANG Xing-xing, LI Yan, YANG Yu-ling and LIU Zhong

College of Systems Engineering, National University of Defense Technology, Changsha 410073, China

Abstract The data form of network can express the entity and the relation between entity and entity. Network structure is common in the real world. It is great significance to study the relationship between nodes and edges in networks. Network representation technology transforms the structure information of network into node vector, which can reduce the complexity of graph representation, and can be effectively applied to tasks such as classification, network reconstruction and link prediction. The SDNE (structural deep network embedding) algorithm proposed in recent years has made outstanding achievements in the field of graph auto-encoder. In view of the limitations of SDNE in weighted and directed networks, this paper proposed a new network representation model based on graph auto-encoder from the perspectives of network structure and measurement index. The concepts of receiving and sending vector are introduced to optimize the decoding part of the neural network, which reduce the parameters of the network to speed up the convergence speed. This paper proposed a measurement index based on the node degree, and reflected the weighted characteristics of the network in the results of the network representation. Experiments on three directed weighted datasets show that the proposed method can achieve better results than the traditional method and the original SDNE method in network reconstruction and link prediction tasks.

Keywords Complex network, Network representation, Network reconstruction, Link prediction, Auto-encoder

1 引言

网络化的数据形式能够自然地表示实体以及实体与实体之间的联系,在日常生活中随处可见^[1-2]。在 Facebook、Twitter 和微博等社交平台中,人与人形成了社交网络,城市间的交通运输构成了交通网络^[3],公司中的员工业务通信构成了企业邮件网络^[4]。基于复杂信息网络的广泛存在,对这类网络信息进行研究与分析具有非常高的学术价值和潜在应用价值。

网络分析中的许多重要任务都涉及对图中的节点或边进

行预测,这就要求使用有效的算法从网络中提取有意义的样式并构造网络特征。近年来,基于机器学习和深度神经网络的方法层出不穷,也取得了不错的成绩。网络表示(图嵌入)方法^[5-7]将网络结构信息转化为节点的低维度向量表示,精确地捕获节点之间的结构关系。研究表明,在节点分类、链路预测和网络重构等多种学习任务中,网络表示方法都取得了优异的成绩。

静态图领域已经存在很多网络表示的方法:奇异值分解^[8](SVD)对网络邻接矩阵进行拉普拉斯变换降维并提取特

收稿日期:2019-06-26 返修日期:2019-10-16 本文已加入开放科学计划(OSID),请扫描上方二维码获取补充信息。

基金项目:国家自然科学基金(61201328,71471175)

This work was supported by the National Natural Science Foundation of China (61201328,71471175).

通信作者:程光权(yang_ma_cn@163.com)

征;Node2vec^[9]通过网络的局部随机游走,将词表示的经典方法 word2vec^[10]运用到网络表示领域;最近提出的 SDNE^[11]算法利用深层自编码对网络进行建模,在保留网络结构信息的基础上,尽可能保留局部的非线性化特征,在网络重构、链路预测等领域中被证实比 Node2vec, SVD 等方法效果更佳^[11]。

SDNE 算法虽然保留了网络中的结构信息,但是对方向信息的处理存在缺陷,算法表示的节点向量丢失了原始网络中与方向相关的关键信息。本文针对 SDNE 在有向网络中的缺陷,从算法结构和计算指标两方面改进了 SDNE 算法的流程,在实验中取得了比原算法更好的实验结果,更好地保留了有向加权网络中的网络特征。

2 算法描述

本节首先介绍网络表示的数学化定义,然后简单介绍 SDNE 算法的基本结构,之后从网络结构和评价指标两个方面给出有向加权网络中的 SDNE 改进方法。

2.1 网络表示定义

定义网络 $G=(V, E)$ 由节点集 $V=\{v_1, \dots, v_n\}$ 和边集 $E=\{e_{i,j}\}_{i,j=1}^n$ 组成, V 中任意两点之间的关联关系记录在边集 E 中。邻接矩阵 A 表示 G 中节点间的连边关系,其 i 行 j 列的 $a_{i,j}$ 值为非负值,表示节点 i 到节点 j 的边的权重值。由于网络为有向有权网络,因此 $a_{i,j}$ 与 $a_{j,i}$ 代表不同的含义,若 $a_{i,j}=0$,表示从节点 i 到节点 j 没有边;若 $a_{i,j}>0$,表示从节点 i 到节点 j 有边,且其值越大表明关系越重要(社交网络)或联系越频繁(通信网络)。 \mathbf{a}_i 表示邻接矩阵的第 i 行,代表节点 i 到网络中所有节点的出度; \mathbf{a}_i' 表示邻接矩阵的第 i 列,代表网络中所有节点到节点 i 的入度。

网络表示旨在将网络中的节点向量由 n 维(或 $2n$ 维)映射到更低的维度空间(k 维, $k=n$),使得到的低维度节点向量能保留节点的结构特征,并能运用到节点分类、链路预测和网络重构等多种网络应用领域,如图 1 所示。

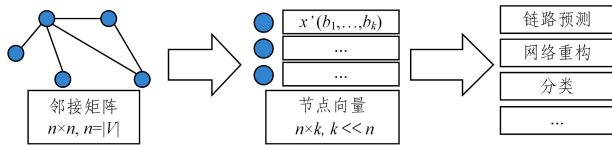


图 1 网络表示的概念图

Fig. 1 Concept map of network embedding

2.2 SDNE 模型概述

SDNE 算法提出了一个半监督的神经网络模型,它具有多层非线性函数,能够捕获到高度非线性的网络结构。该模型以每个节点的邻接矩阵为输入,整个模型可以分为两个部分:一个是由 Laplace 矩阵监督的对第一级相似度进行建模的模块;另一个是由无监督的深层自编码器对第二级相似度关系进行建模的模块。最终,SDNE 算法将深层自编码器的中间层作为节点的网络表示,SDNE 模型如图 2 所示。

Wang 等^[11]提出 SDNE 算法的损失函数由一阶和二阶邻近关系两部分组成,一阶邻近关系使用监督学习来保留网络的局部结构,二阶邻近关系被无监督学习用来捕获全局的网络结构。通过在半监督的神经网络中联合优化它们,可以使模型同时保留局部和全局网络结构特征,从而使算法在稀疏网络中具有较强的鲁棒性。

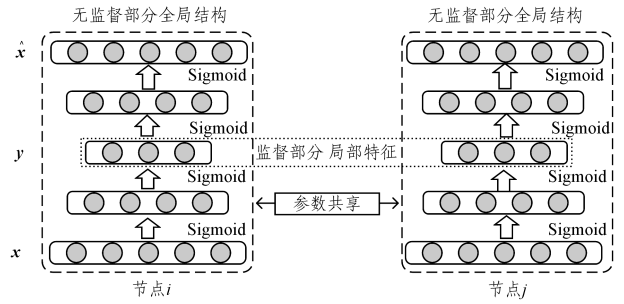


图 2 SDNE 模型的框架

Fig. 2 Structure of SDNE model

一阶邻近关系衡量节点的局部特征,与深层编码器的编码有关,若 $a_{i,j}>0$,即节点 i 和节点 j 之间有连边,则节点 i 和节点 j 经过深层编码器得到的编码应该有相似性。令网络输入向量为 \mathbf{x} ,编码层编码为 \mathbf{y} ,则:

$$loss_{1st} = \sum_{i,j=1}^n a_{i,j} \|\mathbf{y}_i - \mathbf{y}_j\|^2$$

二阶邻近关系衡量节点的全局结构特性与解码器有关,每个节点通过解码器得到的矩阵 $\hat{\mathbf{x}}$ 应该与输入向量 \mathbf{x} 相似,从而保证编码器的编码效果,即:

$$loss_{2nd} = \sum_i \|\mathbf{x}_i - \hat{\mathbf{x}}_i\| \odot \mathbf{b}_i\|^2$$

其中, \odot 表示 Hadamard 积; $\mathbf{b}_i = \{b_{i,j}\}_{j=1}^n$, 若 $a_{i,j}=0$, 则 $b_{i,j}=1$, 若 $a_{i,j}>0$, 则 $b_{i,j}=\theta>1$ 。输入向量 \mathbf{x} 与输出向量 $\hat{\mathbf{x}}$ 的差值与 \mathbf{b}_i 的 Hadamard 积,使损失函数更加关注网络中原有连边的还原度,从而提高了原有连边在稀疏网络中的重要度。

综上,SDNE 的损失函数为:

$$loss = \alpha loss_{1st} + \beta loss_{2nd} + loss_{reg}$$

其中, $loss_{reg}$ 为神经网络参数的 L2 归一化范数:

$$loss_{reg} = \sum_{k=1}^K (\|\mathbf{W}^k\|^2 + \|\hat{\mathbf{W}}^k\|^2) / 2$$

其中, K 是深度神经网络编码器层数, \mathbf{W}^k 表示第 k 层编码器参数, $\hat{\mathbf{W}}^k$ 表示第 k 层解码器参数。

2.3 有向加权网络中网络结构的改进

在无向网络中,任意一个节点 i 到其他节点的出度与入度是相同的,即 $\mathbf{a}_i' = \mathbf{a}_i$;但在有向网络中,节点 i 到其他节点的出度与入度很可能不同。以业务通信网络为例,指挥节点可能会频繁向下级发放相关任务,而下级节点只有在遇到特定情况时,才会向上级反馈,此时网络的节点出度与入度会有明显差距,若仍然用 \mathbf{a}_i 作为深度编码器的输入,则会损失原网络中一半的信息。因此,在有向网络中,本文尝试将节点的出度(邻接矩阵的行)和节点的入度(邻接矩阵的列)拼接在一起,作为深度编码器的输入 $\mathbf{x}_i = [\mathbf{a}_i, \mathbf{a}_i']$ 。

在有向网络中,最重要的是信息的方向。即便将网络的出度矩阵和入度矩阵均考虑在内,传统的解码结构仍无法还原节点向量中蕴含的方向性关系。与 Semi 等^[12]在 Node2vec 中采用的有向向量表示类似,本文修改了传统自编码器的解码网络结构,在通过编码器部分得到节点向量的基础上,将节点向量 \mathbf{y} 分别经过两个相同规格的全连接网络,对每个节点计算其接收向量 \mathbf{x}_{in} 和发出向量 \mathbf{x}_{out} ,并用两向量的向量积 $\mathbf{x}_{out} \times \mathbf{x}_{in}^T$ 作为解码向量 $\hat{\mathbf{x}}$ 。该操作不仅简化了解码器结构,从而加

速神经网络训练,还能更好地提取出同一节点在不同身份(发出和接收)下的节点向量。

SDNE 采用 sigmoid 作为激活函数,但该激活函数两端容易饱和,在神经网络传播过程中节点间区分度较低,且由于加权网络中节点的活跃度差异较大,该现象尤为突出。本文采用 ReLU 激活函数代替 sigmoid,由于其具有分段线性性质,因此前传、后传、求导都是分段线性,有效加快了收敛并降低了误差。SDNE 采用 SGD(Stochastic Gradient Descent)优化神经网络,但 SGD 本身收敛速度慢,且容易陷入鞍点,学习速率的选择对算法效果影响较大。本文采用 Adam(Adaptive moment estimation)优化器以更好地加速神经网络的收敛。此外,改进模型在深度神经网络的每一层后都进行了批标准化处理,避免了梯度爆炸和梯度消失,并提高了算法效果。改进后的 SDNE 模型如图 3 所示。

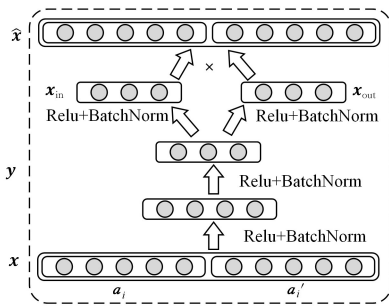


图 3 自编码模型的改进框架

Fig. 3 Improved structure of SDNE model

由于加权有向网络中节点权重差异较大,本文为降低个别高度节点对神经网络优化的影响,在计算二阶邻近关系时,采用绝对值替代平方差,即:

$$loss_{2nd} = \sum_i^n \| (x_i - \hat{x}_i) \odot b_i \|$$

2.4 改进相似度

上一节从深度神经网络结构入手,对 SDNE 算法进行优化,本节从网络表示应用的角度入手,在使用节点向量进行网络重构、链路预测等任务中,在计算节点相似度时更多地强调有权、有向信息,从而改善了算法的实际应用效果。

在有权、有向网络中,根据节点的出度和入度的相对规模,大致可以将节点分为 4 类(见图 4):出度大、入度大的双向活跃节点,出度小、入度大的活跃接收节点,出度大、入度小的活跃发送节点,出度小、入度小的不活跃节点。

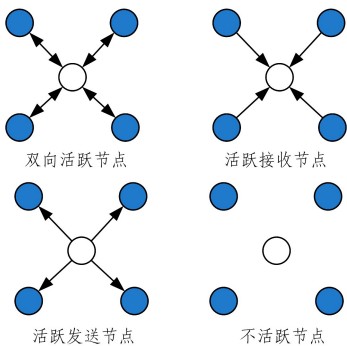


图 4 根据节点度对节点活跃度的划分

Fig. 4 Division of node activity according to node degree

令节点出度为 $I_i = \sum_{j=1}^n a_{i,j}$, 节点入度为 $O_i = \sum_{j=1}^n a_{j,i}$, 对于任意一个节点相似度指标 sim , 本文定义新的加权、有向相似度指标 sim' , 使得:

$$sim'_{i,j} = sim_{i,j} \sqrt{(O_i+1)(I_j+1)}$$

其中, $sim'_{i,j}$ 表示节点 i 到节点 j 的相似度。为了防止节点出度或者入度为 0, 导致原有相似度丢失, 计算时在节点本身出度和入度的基础上加 1, 使得 $\sqrt{(O_i+1)(I_j+1)} > 0$ 。据此, 若节点 i 的出度越大且节点 j 的入度越大, 则节点 i 到节点 j 存在连边的可能性越高。由于计算节点的入度和出度的复杂度仅为 $O(n)$, 在所有与方向和权重相关的应用中, 均可以通过该变换强化网络有向性的影响。

3 实验分析

本节在多个通信数据集中对原 SDNE 算法和改进 SDNE 算法进行实验, 并比较传统衡量指标和改进指标在网络重构和链路预测上的效果。

3.1 实验数据介绍

本文选取安然(Enron)邮件网络^[4]作为实验数据集, 该网络包含了 1999—2003 年间公司人员间的 1148072 个邮件往来。由于邮件允许以自身为收件人, 网络中存在自环。由于时间跨度大且公司人员变动频繁, 本文选取其中 40 个月的数据并分为 3 个独立数据集进行实验。为突出网络主干部分, 本文去除网络自环, 并提取网络中联系密切(节点度大于 400)的部分节点。网络属性如表 1 所列, 网络基本结构如图 5 所示。

表 1 安然邮件网络的基本属性

Table 1 Basic attributes of Enron email networks

编号	节点数	边数	平均度	平均加权重
1	864	93203	17.054	107.874
2	1043	182839	24.712	175.301
3	683	26437	4.765	17.768

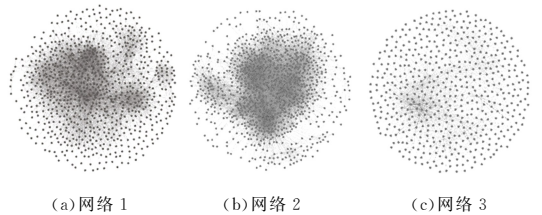


图 5 安然邮件网络图

Fig. 5 Graphs of Enron email networks

3.2 实验参数配置

本次实验搭建 3 层编码器。隐藏层单元数为 500 和 300, 编码层编码为 100 维, 实验采用 ReLU 激活函数, 在每层激活函数后进行批标准化(batchnorm)操作, 批大小 $batchsize=128$, 优化器采用 adam, 学习速率为 $lr=0.01$, 学习代数 $epoch=50$ 。网络 1 中, 一阶邻近参数 $\alpha=0.06$, 二阶邻近参数 $\beta=0.029, \theta=3$; 网络 2 中, 一阶邻近参数 $\alpha=0.12$, 二阶邻近参数 $\beta=0.089, \theta=8$; 网络 3 中, 一阶邻近参数 $\alpha=0.38$, 二阶邻近参数 $\beta=0.09, \theta=6$ 。算法损失函数如图 6 所示。

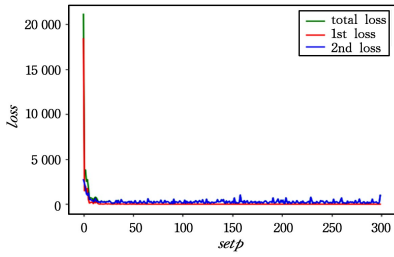


图 6 改进自编码的损失函数图

Fig. 6 Loss function graph of improved SDNE

由图 6 可知,虽然算法的损失函数由两部分组成,但计算初期的损失主要是一阶邻近(1st loss),即局部信息差异,后期误差的主要部分是二阶邻近(2nd loss),即网络全局特征。

3.3 实验结果分析

在对网络表示方法在实际应用中的性能进行评价之前,首先应对该方法在网络重构能力方面进行基本评价,一个好的网络表示方法应该能确保所学到的节点向量能够保持原有的网络结构。由于网络是有序、有向的,在网络重构任务中,我们希望权重更大(更重要)的边被准确地重构出来。本文首先通过网络表示算法得到节点向量,然后计算网络中所有节点的相似性,最后与原始网络中的边权重排序进行比较。本文采取 k-precision^[11] 指标,衡量权重排名前 k 条的边中被准确重构的比例,其在网络 1、网络 2 和网络 3 中的表现如图 7、表 2—表 4 所示。

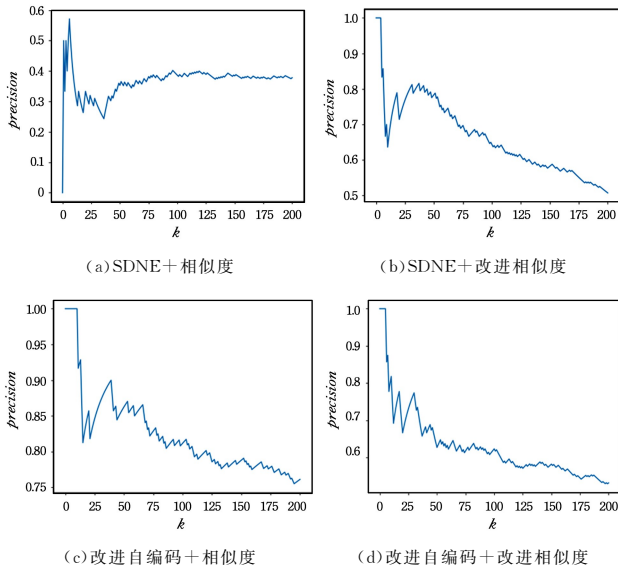


图 7 数据集 1 上多种算法组合的 k-precision(网络重构)

Fig. 7 k-precision indices of algorithm combinations on dataset1 (network reconstruction)

表 2 数据集 1 上各个组合算法的 k-precision(网络重构)

Table 2 k-precision indices of algorithm combinations on dataset1 (network reconstruction)

	10	50	100	200
SDNE	0.4	0.36	0.39	0.38
SDNE+sim	0.7	0.78	0.65	0.51
改进自编码	0.8	0.64	0.62	0.53
改进自编码+sim	1	0.86	0.81	0.76

表 3 数据集 2 中多种算法组合的 k-precision(网络重构)

Table 3 k-precision indices of algorithm combinations on dataset2 (network reconstruction)

	10	50	100	200
SDNE	0.2	0.38	0.35	0.34
SDNE+sim	0.9	0.78	0.71	0.56
改进自编码	0.7	0.68	0.63	0.59
改进自编码+sim	0.9	0.82	0.75	0.66

表 4 数据集 3 中多种算法组合的 k-precision(网络重构)

Table 4 k-precision indices of algorithm combinations on dataset3 (network reconstruction)

	10	50	100	200
SDNE	0.3	0.3	0.25	0.18
SDNE+sim	0.8	0.76	0.64	0.47
改进自编码	0.2	0.42	0.51	0.44
改进自编码+sim	0.9	0.8	0.76	0.71

通过分析可知,在网络重构应用中,单独采用改进自编码结构和改进相似度指标均能对网络重构效果有很大提升,但是同时采取两种措施,k-precision 指标提升最为明显。在数据集 2 和数据集 3 中,单独使用改进相似度指标对 10-precision 的提升尤其明显,说明在单独使用改进相似度指标时,能够很大程度地提高排名靠前节点的预测准确度。在 3 个测试集中,改进自编码+改进相似度组合的 k-precision 效果最好。

在链路预测应用中,选取 80% 的网络数据作为训练集,其余数据作为测试集,通过训练集数据学习节点的向量表示,通过计算节点间的相似度给出边预测排序,并检验排序前 k 的连边在测试集中出现的比例。网络 1、网络 2 和网络 3 中的链路预测 k-precision 指标如图 8、表 5—表 7 所示。

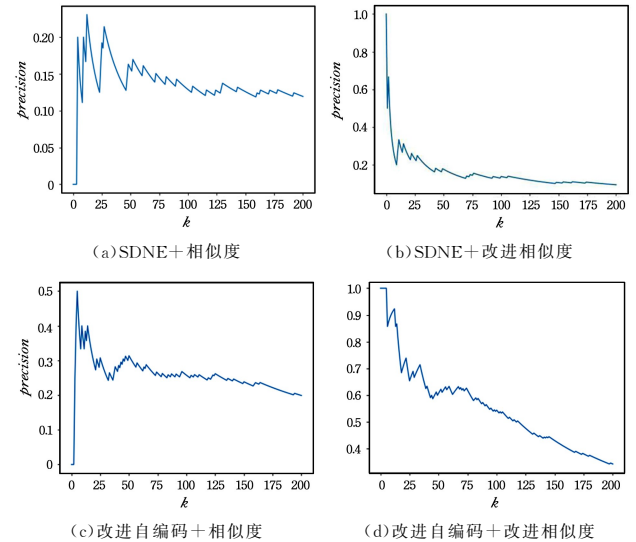


图 8 数据集 1 中多种算法组合的 k-precision(链路预测)

Fig. 8 k-precision indices of algorithm combinations on dataset1 (link prediction)

表 5 数据集 1 中多种算法组合的 k-precision(链路预测)

Table 5 k-precision indices of algorithm combinations on dataset1 (link prediction)

	10	50	100	200
SDNE	0.2	0.16	0.13	0.12
SDNE+sim	0.2	0.18	0.13	0.1
改进自编码	0.4	0.3	0.26	0.2
改进自编码+sim	0.9	0.60	0.54	0.34

表 6 数据集 2 中多种算法组合的 k-precison(链路预测)

Table 6 k-precison indices of algorithm combinations on dataset2
(link prediction)

	10	50	100	200
SDNE	0	0.12	0.08	0.09
SDNE+sim	0.6	0.2	0.18	0.14
改进自编码	0	0.18	0.19	0.13
改进自编码+sim	0.6	0.34	0.29	0.23

表 7 数据集 3 中多种算法组合的 k-precison(链路预测)

Table 7 k-precison indices of algorithm combinations on dataset3
(link prediction)

	10	50	100	200
SDNE	0	0.08	0.06	0.04
SDNE+sim	0.5	0.36	0.26	0.16
改进自编码	0.4	0.18	0.2	0.18
改进自编码+sim	1	0.58	0.4	0.24

通过分析可知,在链路预测应用中,单独采取改进自编码结构和改进相似度指标均能对链路预测效果有很大提升,但是同时采取两种措施,指标提升最为明显。在链路预测指标中,单独采用改进相似度对 10-precision 提升效果非常明显,这与网络重构结果相同。在 3 个测试集中,改进自编码+改进相似度组合的 k-precision 效果最好。

结束语 本文针对网络表示算法 SDNE 在有权、有向网络中的缺陷,从网络结构和衡量指标两个角度入手,提出了改进的 SDNE 模型和基于节点出度、出度的改进衡量指标。经过实验验证,单独运用改进自编码模型和衡量指标,均能在网络重构和链路预测任务中取得更好的效果,联合运用改进策略获得的效果最佳。

所提算法在后续工作中仍有两方面需要继续优化:一方面,该算法虽然更好地利用了网络中的有向、有权信息,但是对节点其他属性信息没有充分利用;另一方面,在动态演化中,每个时间片的网络参数需要进行重复学习,耗费计算时间较多,后续会结合 LSTM^[14]对算法进行优化。

参 考 文 献

- [1] LV L Y. Link Prediction on Complex Networks[J]. Journal of University of Electronic Science and Technology of China, 2010, 39(5): 651-661.
- [2] LV L, ZHOU T. Link prediction in complex networks: A survey [J]. Physica A Statistical Mechanics & Its Applications, 2010, 390(6): 1150-1170.
- [3] MA Y, LIANG X, HUANG J, et al. Intercity Transportation Construction Based on Link Prediction[C]// 2017 IEEE 29th International Conference on Tools with Artificial Intelligence (IC-TAI). IEEE, 2017.
- [4] KLIMT B. The Enron corpus: A new dataset for email classification research[C]// Proc. 15th European Conf. Machine Learning, 2004.
- [5] TU C C, YANG C, LIU Z Y, et al. Network representation

learning: an overview[J]. Scientia Sinical Informationis, 2017 (8): 32-48.

- [6] HAMILTON W L, YING R, LESKOVEC J. Representation Learning on Graphs: Methods and Applications[J]. arxiv:1709.05584.
- [7] GOYAL P, FERRARA E. Graph Embedding Techniques, Applications, and Performance: A Survey[J]. arXiv:1705.02801.
- [8] BELKIN M. Laplacian eigenmaps and spectral techniques for embedding and clustering[J]. Advances in neural information processing systems, 2002, 14(6): 585-591.
- [9] ADITYA GROVER J L. node2vec: Scalable Feature Learning for Networks[J]. arXiv:1607.00653.
- [10] GOLDBERG Y. A Primer on Neural Network Models for Natural Language Processing[J]. arXiv:1510.00726, 2015.
- [11] WANG D, PENG C, ZHU W. Structural Deep Network Embedding[C]// Acm Sigkdd International Conference on Knowledge Discovery & Data Mining. 2016.
- [12] SAMI A, PEROZZI B, AL-RFOU R. Learning Edge Representations via Low-Rank Asymmetric Projections[C]// Proceedings of the 2017 ACM on Conference on Information and Knowledge Management-(CIKM 17). 2017: 1787-1796.
- [13] GOYAL P, KAMRA N, HE X, et al. DynGEM: Deep Embedding Method for Dynamic Graphs[J]. arXiv:1805.11273.
- [14] SUTSKEVER I, VINYALS O, LE Q V. Sequence to Sequence Learning with Neural Networks[J]. arXiv:1511.04868.
- [15] KIM J, PARK H, LEE J E, et al. SIDE: Representation Learning in Signed Directed Networks[C]// the 2018 World Wide Web Conference. 2018.
- [16] CHEN H, PEROZZI B, AL-RFOU R, et al. A Tutorial on Network Embeddings[J]. arXiv:1808.02590.
- [17] CUNCHAO T, XIANGKAI Z, HAO W, et al. A Unified Framework for Community Detection and Network Representation Learning[J]. IEEE Transactions on Knowledge and Data Engineering, 2019, 31(6): 1051-1065.
- [18] GOYAL P, HOSSEINMARDI H, FERRARA E, et al. Capturing Edge Attributes via Network Embedding[J]. IEEE Transactions on Computational Social Systems, 2018, 5(4): 907-917.



MA Yang, born in 1993, postgraduate. His main research interests include link prediction and graph neural networks.



CHENG Guang-quan, born in 1982, Ph.D, is a member of China Computer Federation (CCF). His main research interests include network analysis and machine learning.