

面向推荐系统的图卷积网络^{*}

葛 尧, 陈松灿

(南京航空航天大学 计算机科学与技术学院, 江苏 南京 211106)

通讯作者: 陈松灿, E-mail: s.chen@nuaa.edu.cn



摘 要: 图卷积网络是一种针对图信号的深度学习模型, 由于具有强大的特征表征能力得到了广泛应用. 推荐系统可视为图信号的链接预测问题, 因此近年来提出了使用图卷积网络解决推荐问题的方法. 推荐系统中存在用户与商品间的异质顶点交互和用户(或商品)内部的同质顶点交互, 然而, 现有方法中的图卷积操作要么仅在异质顶点间进行, 要么仅在同质顶点间进行, 留下了提升此类推荐系统性能的空间. 考虑到这一问题, 提出了一种新的基于图卷积网络的推荐算法, 使用两组图卷积操作同时利用两种不同的交互信息, 其中异质顶点卷积用于挖掘交互图谱域中存在的连接信息, 同质顶点卷积用于使相似顶点具有相近表示. 实验结果表明, 该算法比现有算法具有更优的精度.

关键词: 图卷积网络; 图信号; 几何深度学习; 神经网络; 推荐系统

中图法分类号: TP181

中文引用格式: 葛尧, 陈松灿. 面向推荐系统的图卷积网络. 软件学报, 2020, 31(4): 1101–1112. <http://www.jos.org.cn/1000-9825/5928.htm>

英文引用格式: Ge Y, Chen SC. Graph convolutional network for recommender systems. Ruan Jian Xue Bao/Journal of Software, 2020, 31(4): 1101–1112 (in Chinese). <http://www.jos.org.cn/1000-9825/5928.htm>

Graph Convolutional Network for Recommender Systems

GE Yao, CHEN Song-Can

(College of Computer Science and Technology, Nanjing University of Aeronautics and Astronautics, Nanjing 211106, China)

Abstract: Graph convolutional network (GCN) is a deep learning model for graph signal processing and has been used in many real-world applications due to its powerful ability of feature extraction. As the recommendation problem can be viewed as link prediction of graph signals, recently several GCN based methods have been proposed for recommender systems. A recommender system involves two kinds of interactions, with one representing interactions between users and items and the other representing interactions among users (or items). However, existing methods focus on either heterogeneous or homogeneous interactions only, thus their modeling expressiveness is limited. In this study, a new GCN based recommendation algorithm is proposed to jointly utilize these two types of interactions. Specifically, a heterogeneous convolutional operator is used to mine information from the spectrum of user-item graphs, while a homogeneous convolutional operator is used to enforce similar vertices to be similar in the hidden space. Finally, the experiments on benchmark datasets show that the proposed method achieves better performance compared with several state-of-the-art methods.

Key words: graph convolutional network; graph signal; geometric deep learning; neural network; recommender system

卷积神经网络(convolutional neural network, 简称 CNN)^[1]具有强大的特征表征能力, 因而在诸如计算机视觉、自然语言处理等领域得到了极大关注^[2,3]. 文本、图像和视频均是定义在规则网格(regular grid)上的数据, 它们能对应地视为分布在一维、二维和三维网格支撑集上. CNN 能够方便地运算可归因于这些网格的规则性.

^{*} 基金项目: 国家自然科学基金(61672281, 61732006)

Foundation item: National Natural Science Foundation of China (61672281, 61732006)

本文由“非经典条件下的机器学习方法”专题特约编辑高新波教授、黎铭教授、李天瑞教授推荐.

收稿时间: 2019-05-31; 修改时间: 2019-07-29; 采用时间: 2019-09-20; jos 在线出版时间: 2020-01-10

CNKI 网络优先出版: 2020-01-14 09:53:29, <http://kns.cnki.net/kcms/detail/11.2560.TP.20200114.0953.013.html>

然而,现实中,除这些规则网格数据之外,还存在一类重要的称为图信号(graph signal)^[4]的数据(下称图信号),它分布或定义在不规则网格(irregular grid)支撑集上.针对图信号的深度学习又称为几何深度学习(geometric deep learning)^[5].一方面,图信号可视为图及其顶点上的信号集合;另一方面,图信号也可视为一组非独立同分布的数据点,是一类非传统型数据,它们之间的关系用图的链接表示.如何将传统 CNN 推广到能够处理更复杂的图信号的卷积网络,即图卷积网络(graph convolutional network,简称 GCN),利用其强大的特征表征能力提升学习效果,正得到越来越多的关注.

目前已有许多工作围绕 GCN 展开研究,并在理论和应用上取得了丰富成果^[6-16].将一个问题用图信号刻画后,就可根据问题特点设计相应的 GCN 来解决.学者网络^[8]、社交网络^[11]、分子活性预测^[12]和推荐系统^[13]都是 GCN 的典型应用场景.具体到推荐系统,涉及的用户和商品可视为图顶点,用户对商品的评分可视为边(链接),而用户和商品特征视为分布在顶点集(图)上的信号,如此,推荐问题便转化为图的链接预测问题.值得注意的是,推荐系统中存在两种图:一种是异质顶点交互图,反映用户对商品的行为,例如评分、购买等;另一种是同质顶点交互图,反映用户(商品)间的相似性,例如朋友关系、商品特征相似等.两种图都包含部分信息,兼顾两者以实现信息互助对推荐系统至关重要.然而,现有的基于 GCN 的推荐系统要么仅关注异质图,要么仅关注同质图,缺少联合利用两者的统一框架,由此为我们留下了深入利用图提升推荐性能的空间.因此,本文的目的就是提出能够统一利用这两种图的 GCN,通过它们的互惠互利提升推荐效果.下面我们首先介绍现有的两类方法.

异质顶点交互 GCN 类方法(hetero-GCN)在异质顶点间进行图卷积操作.将 m 个用户对 n 个商品的评分视为一个 $(m+n) \times (m+n)$ 的二部图,此类方法重点挖掘交互图谱域中蕴含的连接信息,使用 GCN 直接从二部图中提取特征.GC-MC^[17]和 SpectralCF^[18]是两种代表方法,两者都以顶点特征为信号,以评分信息为图进行图卷积操作.此类方法仅使用异质顶点交互信息,忽略了顶点相似性信息,在评分过少时会遇到冷启动(cold start)问题.

同质顶点交互 GCN 类方法(homo-GCN)在同质顶点间进行图卷积操作.将 m 个用户对 n 个商品的评分视为一个 $m \times n$ 的矩阵,使用评分或特征信息构建 $m \times m$ 的行图(row graph)和 $n \times n$ 的列图(column graph),分别代表用户和商品相似度.此类方法认为相似用户(商品)的表示向量应当相近,因此在相似的同质顶点间进行图卷积,使其信号平滑.RGCNN^[19]和 GCMC-BEP^[20]是两种代表方法,两者都从评分矩阵获取初始特征,分别在行图和列图上进行图卷积获取新特征表示.此类方法将评分信息视为提供特征的矩阵,没有将其视为图,使评分图中蕴含的连接信息未能得到利用.

本文第 1 节介绍 Hetero-GCN 框架及代表方法.第 2 节介绍 Homo-GCN 框架及代表方法.第 3 节提出一种联合利用异质与同质交互信息的 GCN 推荐算法.第 4 节在真实数据集上进行实验,验证本文方法优于现有方法.最后总结全文,并对未来工作进行展望.

1 Hetero-GCN

1.1 模型框架

设 m 个用户对 n 件商品的评分矩阵为 $R \in R^{m \times n}$, 评分取值范围为 $\{1, 2, \dots, L\}$, 用户和商品特征为 $X'_u \in R^{m \times d_u}$ 和 $X'_i \in R^{n \times d_i}$. 为每一级评分构建用户-商品交互图:

$$G_{hetero}^l = \begin{bmatrix} & R^l \\ (R^l)^T & \end{bmatrix}, l = 1, 2, \dots, L \quad (1)$$

其中, $G_{hetero}^l \in R^{(m+n) \times (m+n)}$, 且 $R_{i,j}^l = \begin{cases} 1, & R_{i,j} = l \\ 0, & R_{i,j} \neq l \end{cases}$.

记 G_{hetero} 为交互图集合: $G_{hetero} = \{G_{hetero}^l | l = 1, 2, \dots, L\}$. 为简化记号, 记 $X_u = [X'_u, 0] \in R^{m \times d}$, $X_i = [0, X'_i] \in R^{n \times d}$, 则 $X = [X_u^T \ X_i^T]^T$ 代表所有用户和商品特征. 如此, 推荐系统可用图信号 $\{G_{hetero}, X\}$ 表示, 其中, G_{hetero} 代表图, X 代表顶点信号, 推荐问题转化为对 G_{hetero} 的链接预测问题. 用 GCN 为顶点学习向量表示(又称嵌入向量), 利用嵌入向量进行链接预测.

将待学习嵌入向量记为 $\mathbf{Z} = [\mathbf{Z}_u^T \quad \mathbf{Z}_i^T]^T \in R^{(m+n) \times c}$, 则 Hetero-GCN 中卷积操作为

$$\mathbf{Z} = \text{Conv}(\mathbf{G}_{\text{hetero}}, \mathbf{X}) \quad (2)$$

图 1 给出了一个 Hetero-GCN 的卷积操作示例,图中左半部分表示用户商品评分图,右半部分表示图卷积操作过程.图卷积在异质顶点间进行,用户 u_1 的新特征表示来源于其评分的商品 i_1 、 i_2 、 i_4 ;商品 i_4 的新特征来源于为其评分的用户 u_1 和 u_2 .

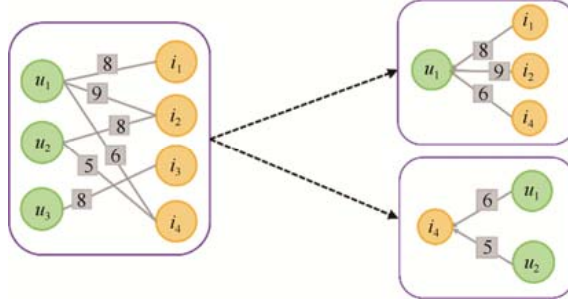


Fig.1 Convolution operator in hetero-GCN

图 1 Hetero-GCN 中的卷积操作

1.2 代表方法

1.2.1 GC-MC

设有图信号 $\{\mathbf{G}, \mathbf{x}\}$, \mathbf{G} 的图拉普拉斯矩阵为 $\mathbf{L} = \mathbf{I} - \mathbf{D}^{-\frac{1}{2}} \mathbf{G} \mathbf{D}^{-\frac{1}{2}} = \mathbf{U} \mathbf{A} \mathbf{U}^T$, 图傅里叶变换定义为 $\hat{\mathbf{x}} = \mathbf{U}^T \mathbf{x}$ [21]. 根据卷积定理,两个信号卷积的傅里叶变换等于它们傅里叶变换的乘积,滤波器 \mathbf{g}_θ 与信号 \mathbf{x} 卷积的结果为 [6]

$$\mathbf{z} = \mathbf{g}_\theta * \mathbf{x} = \mathbf{U} \mathbf{g}_\theta(\mathbf{\Lambda}) \mathbf{U}^T \mathbf{x} \quad (3)$$

式(3)的图卷积操作存在 3 个问题:需要特征值分解、滤波器参数数量并非常数以及滤波器未被限定在局部.文献[7]使用 Chebyshev 多项式展开 $\mathbf{g}_\theta(\mathbf{\Lambda})$ 并限定到 k 阶解决了上述问题.文献[8]进一步将多项式限定到一阶简化计算:

$$\mathbf{z} = \theta \left(\mathbf{I} + \mathbf{D}^{-\frac{1}{2}} \mathbf{G} \mathbf{D}^{-\frac{1}{2}} \right) \mathbf{x} \quad (4)$$

将信号扩展到多通道并添加非线性变换,使用重归一化技巧加强数值稳定性,得到一阶近似图卷积:

$$\mathbf{Z} = \sigma \left(\tilde{\mathbf{D}}^{-\frac{1}{2}} \tilde{\mathbf{G}} \tilde{\mathbf{D}}^{-\frac{1}{2}} \mathbf{X} \mathbf{W} \right) \quad (5)$$

其中, $\tilde{\mathbf{G}} = \mathbf{G} + \mathbf{I}$, $\tilde{\mathbf{G}}$ 是 $\tilde{\mathbf{D}}$ 的度矩阵, $\mathbf{X} \in R^{N \times d}$ 是输入图信号, $\mathbf{Z} \in R^{N \times c}$ 是输出图信号, $\mathbf{W} \in R^{d \times c}$ 是待学习滤波器参数.

式(5)所示卷积运算有直觉解释:对图顶点 v_i ,将其所有邻接点 v_j 上的特征(信号) \mathbf{x}_j 按边权 g_{ij} 相加,经非线性变换后作为顶点 v_i 的新特征.这与图像处理中经典卷积工作原理相同,故式(5)成为目前使用最多的图卷积定义,GC-MC 卷积即采用此定义:

$$\begin{bmatrix} \mathbf{Z}_u \\ \mathbf{Z}_i \end{bmatrix} = \sigma \left(\text{accum}_l \left[(\mathbf{D}')^{-\frac{1}{2}} \mathbf{M}' (\mathbf{D}')^{-\frac{1}{2}} \begin{bmatrix} \mathbf{X}_u \\ \mathbf{X}_i \end{bmatrix} \right] \mathbf{W} \right) \quad (6)$$

其中, $\mathbf{M}' = \begin{bmatrix} & \mathbf{R}' \\ (\mathbf{R}')^T & \end{bmatrix}$, \mathbf{D}' 是 \mathbf{M}' 的度矩阵. $\text{accum}[\cdot]$ 代表聚合函数,例如 $\text{stack}[\cdot]$ 串联函数,或 $\text{sum}[\cdot]$ 求和函数.

SpectralCF.

SpectralCF 用于解决隐式推荐问题,此类问题中仅有用户对商品浏览、购买等行为信息,没有显式评分.其中, $\mathbf{R} \in \{0, 1\}^{m \times n}$ 代表用户与商品是否存在交互,因此 $\mathbf{G}_{\text{hetero}}$ 仅包含一个交互图.

SpectralCF 将式(3)中的 $g_\theta(\Lambda)$ 多项式展开并限定到一阶:

$$z = \theta(UU^T + UAU^T)x = \theta(I + G)x \quad (7)$$

将信号扩展到多通道并添加非线性变换得到:

$$Z = \sigma((I + G)XW) \quad (8)$$

SpectralCF 中图卷积操作为

$$\begin{bmatrix} Z_u \\ Z_i \end{bmatrix} = \sigma((I + M) \begin{bmatrix} X_u \\ X_i \end{bmatrix} W) \quad (9)$$

其中, $M = \begin{bmatrix} & R \\ R^T & \end{bmatrix}$.

事实上,式(6)和式(9)图卷积操作的唯一区别为是否在邻接图中添加自环及是否归一化.GC-MC 和 SpectralCF 可视为 Hetero-GCN 在显式和隐式推荐系统上的具体实现.

此外,PinSage^[13]也是一种 Hetero-GCN 模型,其面向隐式推荐问题.PinSage 采用与 GC-MC 类似的卷积定义,但着重于超大规模推荐系统的工业级实现.

2 Homo-GCN

2.1 模型框架

与第 1.1 节相似,设有评分矩阵 R 、用户和商品特征 X'_u, X'_i . 另有用户和商品相似度矩阵 $G_u \in R^{m \times m}$, $G_i \in R^{n \times n}$.

相似度矩阵可通过多种途径获得,例如外部社交网络信息、评分向量相似度和特征相似度等.这里假设已获得 G_u 和 G_i ,则同质相似图为

$$G_{homo} = \begin{bmatrix} G_u & \\ & G_i \end{bmatrix} \quad (10)$$

其中, $G_{homo} \in R^{(m+n) \times (m+n)}$.

与 Hetero-GCN 将 R 视为图不同,Homo-GCN 将 R 视为矩阵.Homo-GCN 从 R 中提取信息作为顶点信号.例如将 R 中每行和每列分别作为用户和商品特征,或对 R 低秩分解后将左右因子作为特征.推荐系统可用图信号 $\{G_{homo}, (X, R)\}$ 表示.

从 X 和 R 中获得顶点信号 X_u 和 X_i ,随后分别在 G_u 和 G_i 上进行图卷积,得到嵌入向量.

$$Z = \text{Conv}(G_{homo}, (X, R)) \quad (11)$$

图 2 给出了一个 Homo-GCN 的卷积操作示例,图中左半部分表示同质交互图和评分矩阵,右半部分表示图卷积操作过程.评分矩阵仅提供初始特征,不以图的形式参与卷积过程.图卷积在同质顶点间进行,用户 u_1 的新特征表示来源于相似用户 u_2 和 u_3 ;商品 i_3 的新特征来源于相似商品 i_1 和 i_2 .

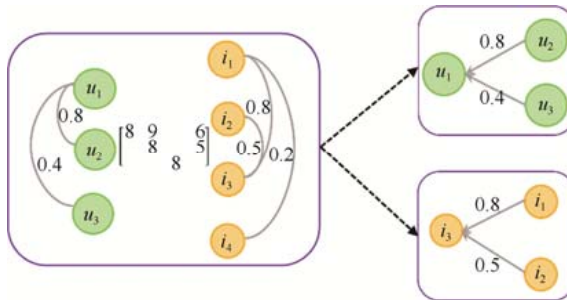


Fig.2 Convolution operator in homo-GCN

图 2 Homo-GCN 中的卷积操作

2.2 代表方法

2.2.1 RGCNN

RGCNN 对评分矩阵 R 低秩分解 $R = X_u X_i^T$, 得到因子矩阵作为用户和商品顶点信号. 此方法使用 Chebyshev 多项式图卷积操作. 在式(3)中, 将 $g_\theta(\Lambda)$ 展开:

$$g_\theta(\Lambda) = \sum_{k=0}^{K-1} \theta_k T_k(\tilde{\Lambda}) \quad (12)$$

随后可借助 Chebyshev 多项式的递归性质简化计算, 具体计算方式见文献[7]. RGCNN 中图卷积操作为

$$\begin{bmatrix} Z_u \\ Z_i \end{bmatrix} = \begin{bmatrix} \text{ChebyNet}(G_u, X_u) \\ \text{ChebyNet}(G_i, X_i) \end{bmatrix} \quad (13)$$

即用户和商品顶点信号分别在 G_u 和 G_i 上图卷积.

2.2.2 GCMC-BEP

GCMC-BEP 同样对评分矩阵 R 低秩分解获取初始特征. 随后在 G_u 和 G_i 上对 X_u 和 X_i 分别进行图卷积获得新特征表示. 图卷积定义采用式(5).

$$\begin{bmatrix} Z_u \\ Z_i \end{bmatrix} = \sigma \left(\begin{bmatrix} \tilde{G}_u \\ \tilde{G}_i \end{bmatrix} \begin{bmatrix} X_u \\ X_i \end{bmatrix} W \right) \quad (14)$$

其中, \tilde{G}_u 和 \tilde{G}_i 分别是对 G_u 和 G_i 添加自环并归一化后的图.

文献[8]指出, 式(14)中使用的一阶图卷积实际上是对式(13)中的 Chebyshev 图卷积的近似, 故 GCMC-BEP 可视为 RGCNN 的简化版本.

此外, GCNCF^[14]也是一种 Homo-GCN 模型, GCNCF 在每一级评分上构建同质交互图, 并进行图卷积操作.

3 GCN4RS

异质顶点图卷积从交互图谱域中提取信息, 同质顶点图卷积使相似顶点有相近表示. 为同时利用异质和同质顶点交互信息, 通过信息互助提高推荐系统性能, 本文提出了 GCN4RS(graph convolutional network for recommender systems)算法. GCN4RS 采用自编码器(autoencoder)^[22]框架, 结构如图 3 所示: 编码器包含提取图信息的 GCN 层和提取特征信息的全连接层, 解码器根据嵌入向量相似度预测链接存在与否.

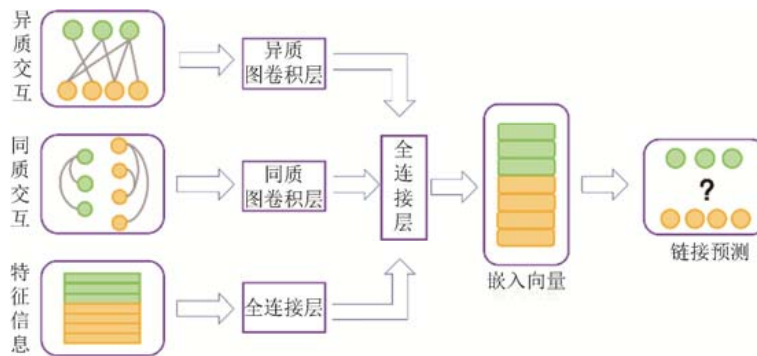


Fig.3 Framework of GCN4RS

图 3 GCN4RS 框架

3.1 编码器

在使用图信号建模推荐系统时, 链接刻画异质顶点交互和同质顶点交互信息, 顶点信号刻画特征信息, 使用 GCN 作为编码器能够统一利用这些信息, 通过它们的互惠互利提升推荐系统性能. 同时, 文献[8]指出, 在没有顶点特征时, 以顶点序号独热编码(one-hot encoding)为顶点信号的 GCN 就可获得颇具竞争力的效果, 文献[17]同

样指出,处理推荐问题时,以独热编码为顶点信号,以顶点特征为单独信息源可获得更优性能,因此本文也采用这种做法.如此,我们的编码器应当统一利用异质交互、同质交互和顶点特征这3种不同的信息.

我们使用图卷积层提取异质交互信息:

$$\mathbf{Z}_{hetero}^l = (\mathbf{D}_{hetero}^l)^{-1} \begin{bmatrix} \mathbf{R}^l \\ \mathbf{R}^{l^T} \end{bmatrix} \begin{bmatrix} \tilde{\mathbf{X}}_u \\ \tilde{\mathbf{X}}_i \end{bmatrix} \mathbf{W}_1 \quad (15)$$

其中, \mathbf{D}_{hetero}^l 是 $\begin{bmatrix} \mathbf{R}^l \\ \mathbf{R}^{l^T} \end{bmatrix}$ 的度矩阵,用于归一化, $\mathbf{R}_{i,j}^l = \begin{cases} 1, & \mathbf{R}_{i,j} = l \\ 0, & \mathbf{R}_{i,j} \neq l \end{cases}$ 代表在 l 评分上用户与商品的交互, $\tilde{\mathbf{X}}_u$ 和 $\tilde{\mathbf{X}}_i$ 是用户和商品独热编码, \mathbf{W}_1 是待学习的参数.异质嵌入向量联合所有评分 l 上的异质交互信息: $\mathbf{Z}_{hetero} = \text{stack}_l[\mathbf{Z}_{hetero}^l]$.

类似地,使用图卷积层提取同质交互信息:

$$\mathbf{Z}_{homo} = (\mathbf{D}_{homo})^{-1} \begin{bmatrix} \mathbf{G}_u \\ \mathbf{G}_i \end{bmatrix} \begin{bmatrix} \tilde{\mathbf{X}}_u \\ \tilde{\mathbf{X}}_i \end{bmatrix} \mathbf{W}_2 \quad (16)$$

其中, \mathbf{D}_{homo} 是 $\begin{bmatrix} \mathbf{G}_u \\ \mathbf{G}_i \end{bmatrix}$ 的度矩阵,用于归一化, \mathbf{G}_u 和 \mathbf{G}_i 分别代表用户交互矩阵和商品交互矩阵, $\tilde{\mathbf{X}}_u$ 和 $\tilde{\mathbf{X}}_i$ 是用户和商品独热编码, \mathbf{W}_2 是待学习的参数.在没有外部同质交互信息输入,如社交网络、商品关联时, \mathbf{G}_u 和 \mathbf{G}_i 需使用已有信息构建,如利用顶点特征相似性、评分模式相似性等.本文实验部分将给出一种供参考的选取方法.

顶点特征由全连接层接入网络:

$$\mathbf{Z}_{feat} = \mathbf{W}_3 \begin{bmatrix} \mathbf{X}_u \\ \mathbf{X}_i \end{bmatrix} + \mathbf{b}_3 \quad (17)$$

最后使用一个全连接层联合利用3种不同的信息 $\mathbf{Z}_{hetero}, \mathbf{Z}_{homo}, \mathbf{Z}_{feat}$:

$$\mathbf{Z} = \sigma(\mathbf{W}_4 [\mathbf{Z}_{hetero}, \mathbf{Z}_{homo}, \mathbf{Z}_{feat}] + \mathbf{b}_4) \quad (18)$$

GCN4RS的异质图卷积层和同质图卷积层均使用单层(1-layer)图卷积.编码器输出 \mathbf{Z} 即顶点嵌入向量.在模型训练时, \mathbf{Z} 被送入解码器中重建输入,并通过误差反向传播不断更新.训练结束后, \mathbf{Z} 中已嵌入推荐系统的交互和特征信息,可用于完成推荐任务.

值得注意的是,虽然 GCN4RS 进行了两次图卷积操作,但输入信号 $\begin{bmatrix} \tilde{\mathbf{X}}_u \\ \tilde{\mathbf{X}}_i \end{bmatrix}$ 是独热编码向量,仅含有 $m+n$ 个非

零元素,且推荐系统中的图多为仅含 $O(m+n)$ 个顶点的稀疏图,卷积操作的核心部分矩阵乘法 $\begin{bmatrix} \mathbf{R} \\ \mathbf{R}^T \end{bmatrix} \begin{bmatrix} \tilde{\mathbf{X}}_u \\ \tilde{\mathbf{X}}_i \end{bmatrix}$ 复杂度仅有 $O(m+n)$,相比于现有方法未增加额外计算开销.

3.2 解码器与损失函数

解码器根据编码器输出的嵌入向量 \mathbf{Z} 进行链接预测,使用用户向量 \mathbf{z}_i^u 和商品向量 \mathbf{z}_j^i 预测评分 \mathbf{R}_{ij} . GCN4RS 使用双线性解码器预测评分 \mathbf{R}_{ij} 为 l 的概率:

$$p(\mathbf{R}_{ij} = l) = \frac{\exp((\mathbf{z}_i^u)^T \mathbf{Q}_l \mathbf{z}_j^i)}{\sum_{s=1}^L \exp((\mathbf{z}_i^u)^T \mathbf{Q}_s \mathbf{z}_j^i)} \quad (19)$$

其中, \mathbf{z}_i^u 和 \mathbf{z}_j^i 分别代表第 i 个用户和第 j 个商品的嵌入向量. $\mathbf{Q}_l \in \mathbb{R}^{c \times c}$ 是可训练的参数,相比于直接使用 \mathbf{z}_i^u 与 \mathbf{z}_j^i 的内积,引入 \mathbf{Q}_l 可增强模型的拟合能力.

解码器输出的预测值 $\tilde{\mathbf{R}}_{ij}$ 为评分关于概率 $p(\tilde{\mathbf{R}}_{ij} = l)$ 的期望:

$$\tilde{\mathbf{R}}_{ij} = \sum_{l=1}^L l \cdot p(\tilde{\mathbf{R}}_{ij} = l) \quad (20)$$

优化目标选用交叉熵损失:

$$\mathcal{L} = \sum_{\Omega_{ij}=1} \sum_{l=1}^L \mathbb{I}[\mathbf{R}_{ij} = l] \cdot \log p(\mathbf{R}_{ij} = l) \quad (21)$$

其中, $\Omega_{ij}=1$ 代表 \mathbf{R}_{ij} 位置的元素已知, $\mathbb{I}[\cdot]$ 代表指示函数. 使用梯度类优化算法训练.

3.3 算法流程

算法 1. GCN4RS.

输入: 评分矩阵 $\mathbf{R} \in R^{m \times n}$, 用户特征 $\mathbf{X}'_u \in R^{m \times d_u}$, 商品特征 $\mathbf{X}'_i \in R^{n \times d_i}$, 嵌入向量维度 c , 迭代次数 epoch.

输出: 用户和商品嵌入向量 $\mathbf{Z} = \begin{bmatrix} \mathbf{Z}'_u & \mathbf{Z}'_i \end{bmatrix}^T \in R^{(m+n) \times c}$.

- 1) 对用户和商品序号独热编码获得顶点信号 $\tilde{\mathbf{X}}_u \in R^{m \times (m+n)}$, $\tilde{\mathbf{X}}_i \in R^{n \times (m+n)}$
- 2) 将用户和商品特征填充到同一维度: $\mathbf{X}_u = [\mathbf{X}'_u, 0] \in R^{m \times d}$, $\mathbf{X}_i = [0, \mathbf{X}'_i] \in R^{n \times d}$
- 3) 构建异质交互图 $\mathbf{G}_{hetero}^l = \begin{bmatrix} & \mathbf{R}^l \\ \mathbf{R}^{l^T} & \end{bmatrix}$, $l=1, 2, \dots, L$. 其中, $\mathbf{R}_{ij}^l = \begin{cases} 1, & \mathbf{R}_{i,j} = l \\ 0, & \mathbf{R}_{i,j} \neq l \end{cases}$
- 4) 构建同质交互图 $\mathbf{G}_{homo} = \begin{bmatrix} \mathbf{G}_u & \\ & \mathbf{G}_i \end{bmatrix}$, 构建方法见第 3.1 节
- 5) for $i=1$: epoch do
- 6) 执行异质和同质图卷积 $\mathbf{Z}'_{hetero} = (\mathbf{D}_{hetero}^l)^{-1} \begin{bmatrix} & \mathbf{R}^l \\ \mathbf{R}^{l^T} & \end{bmatrix} \begin{bmatrix} \tilde{\mathbf{X}}_u \\ \tilde{\mathbf{X}}_i \end{bmatrix} \mathbf{W}_1$, $\mathbf{Z}'_{homo} = (\mathbf{D}_{homo})^{-1} \begin{bmatrix} \mathbf{G}_u & \\ & \mathbf{G}_i \end{bmatrix} \begin{bmatrix} \tilde{\mathbf{X}}_u \\ \tilde{\mathbf{X}}_i \end{bmatrix} \mathbf{W}_2$
- 7) 执行特征的全连接层运算 $\mathbf{Z}_{feat} = \mathbf{W}_3 \begin{bmatrix} \mathbf{X}_u \\ \mathbf{X}_i \end{bmatrix} + b_3$
- 8) 根据 $\mathbf{Z} = \sigma(\mathbf{W}_4 [\mathbf{Z}_{hetero}, \mathbf{Z}_{homo}, \mathbf{Z}_{feat}] + b_4)$ 获得嵌入向量 \mathbf{Z}
- 9) 根据式(21)计算损失并反向传播梯度更新式(19)中的解码器参数 \mathbf{Q}_l 和嵌入向量 \mathbf{Z}
- 10) 反向传播更新编码器参数 $\mathbf{W}_1, \mathbf{W}_2, \mathbf{W}_3, \mathbf{W}_4, \mathbf{W}_5, \mathbf{W}_6$
- 11) 返回 5)
- 12) end for

3.4 GCN类推荐算法对比

3 类方法都用图信号对推荐系统建模, 图信号由图结构和顶点信号组成, 因此, 3 类方法的主要区别在于将推荐系统中的何种信息视为图信号的何种成分. 如表 1 所示, Hetero-GCN^[13,17,18]在评分图上对顶点特征进行卷积, Homo-GCN^[14,19,20]在相似度图上对来自评分矩阵的特征进行卷积. 我们的方法 GCN4RS 将评分和顶点相似性都视为图.

Table 1 Comparison of the way to use information in RS

表 1 对推荐系统中信息的利用方式比较

	Hetero-GCN	Homo-GCN	GCN4RS
评分	图	顶点信号	图
顶点特征	顶点信号	无	顶点信号
顶点相似性	无	图	图

Hetero-GCN 将评分矩阵 \mathbf{R} 视为含有 $m+n$ 个顶点的二部图, 其中链接仅存在于用户和商品顶点间. 相对于矩阵, 图包含更多信息, 例如图谱域中包含对信号高频和低频的分离^[4], 尽可能多地用图表示数据更有利于对信息的挖掘. 用户和商品特征信息则以顶点信号的形式得以利用. 但 Hetero-GCN 未利用顶点相似性信息, 在评分数量较少时会遇到冷启动问题.

Homo-GCN 分别使用含有 m 和 n 个顶点的图刻画用户相似度和商品相似度, 边权代表顶点间的相似程度.

由于图卷积操作会使特征在相邻顶点间流动交互,相似的顶点会有相近的特征表示.评分信息以顶点信号的形式得以利用,例如将评分矩阵的每行和每列分别作为用户和商品顶点的初始信号,或对评分矩阵分解后将因子矩阵作为顶点信号.Homo-GCN 未利用用户和商品的特征信息,一种简单的修改方式是将特征信息也视为顶点信号.此外,图相对于矩阵可刻画更多信息,Homo-GCN 以矩阵而不是图的形式利用评分,丢失了评分图谱域中蕴含的连接信息.

GCN4RS 用两种图分别刻画异质和同质顶点交互信息,利用 GCN 挖掘图中蕴含的信息,使推荐系统中的交互信息得以充分利用.同时,顶点特征信息,例如用户资料、商品描述等则以顶点信号的形式得以利用.在评分数量较少时,顶点相似信息可在一定程度上缓解冷启动问题;在顶点特征和顶点相似信息不足时,评分图谱域中蕴含的连接信息可作为有力补充.如此,异质和同质交互信息在 GCN4RS 框架中实现了互助.本文实验将证明这一信息互助可以提升推荐系统的性能.

3 类算法虽利用了不同信息,但本质上都基于图卷积网络,主要计算代价都来自图卷积运算.与第 3.1 节中分析相同,推荐系统多为边数正比于顶点数的稀疏图.对含有 m 个用户和 n 个商品的推荐系统,Hetero-GCN 在含有 $O(m+n)$ 条边的异质二分图上进行图卷积,Homo-GCN 在含有 $O(m)+O(n)$ 条边的两个同质图上进行图卷积,GCN4RS 同时进行两种图卷积,边的数量为 $O(m+n)+O(m)+O(n)$.可见,GCN4RS 的时间复杂度与 Hetero-GCN 和 Homo-GCN 同阶,都为 $O(m+n)$.虽然需要进行两组图卷积运算,但异质图 and 同质图上的图卷积运算过程完全独立,很容易并行执行,GCN4RS 的运算时间 T_{GCN4RS} 可从 $T_{hetero} + T_{homo}$ 降低到 $\max(T_{hetero}, T_{homo})$.

4 实验与结果

4.1 实验设置

为了验证 GCN4RS 的性能,在 4 个通用推荐系统数据集上进行了实验,数据集的基本信息见表 2.Flixster、Douban 和 YahooMusic 数据集使用文献[19]提供的经过处理的子集,均包含 3 000 用户和 3 000 商品.MovieLens 按 0.8/0.2 划分训练集和测试集,其他 3 个数据集按 0.9/0.1 划分.

Table 2 Experimental datasets

表 2 实验数据

数据集	用户数量	商品数量	评分数量	评分密度	评分级别
Flixster	3 000	3 000	26 173	0.002 9	0.5,1,...,5
Douban	3 000	3 000	136 891	0.015 2	1,2,...,5
YahooMusic	3 000	3 000	5 335	0.000 6	1,2,...,100
MovieLens	943	1 682	100 000	0.063 0	1,2,...,5

异质交互图的构建方法为:为每一级评分构建一个 0-1 邻接图,Flixster、Douban、YahooMusic 和 MovieLens 数据集分别含有 10、5、71、5 级不同的评分(YahooMusic 数据集中仅有 71 种不同的评分出现),故分别构建 10、5、71、5 个 0-1 交互图.

同质交互图 G_u 和 G_i 的构建方法为:在共同评分数多于 K 的顶点间加入链接,权值为特征相似度. K 越小,链接信息越丰富,但计算开销随之增加,需在两者间加以权衡.依评分密度不同,Flixster、Douban、YahooMusic 和 MovieLens 数据集的阈值 K 分别选为 5、15、2、30.

我们的模型使用 TensorFlow 实现.经交叉验证后选用如下超参数设置. Z_{hetero} 和 Z_{homo} 维度设为 200, Z_{feat} 和 Z_{out} 维度设为 64,顶点信号 dropout 概率设为 0.7,激活函数选用 ReLU().使用 Adam 优化器^[23],学习率为 0.01.

对 MovieLens 数据集,顶点特征不进行 dropout,运行 1 000 轮迭代;对其他 3 个数据集,顶点特征 dropout 概率设为 0.7,运行 200 轮迭代.参照文献[17]中的建议,在参数学习过程中加入衰减因子 0.995 的指数移动平均.

对比方法选取异质交互类算法 GC-MC^[17]、同质交互类算法 RGCNN^[19]、GCNCF^[14],并选取矩阵补全 MC^[24]、几何矩阵补全 GMC^[25]、交替最小二乘几何矩阵补全 GRALS^[26]等算法作为参照.此外,还设置两种 GCN4RS 变种进行比对分析,其中,GCN4RS-hetero 仅进行异质顶点卷积,GCN4RS-homo 仅进行同质顶点卷积.

评价指标采用 RMSE:

$$\text{RMSE} = \sqrt{\frac{1}{n} \sum_{i=1}^n (R_{ij} - \tilde{R}_{ij})^2} \quad (22)$$

其中, n 为测试样本数量, R_{ij} 和 \tilde{R}_{ij} 分别代表第 i 个样本的真实和预测评分. 实验中对算法随机初始化并运行 5 次, RMSE 取 5 次的平均值.

4.2 实验结果

4.2.1 参数选取

如第 4.1 节所述, 为使用同质交互信息, 需对 4 个数据集分别构建同质交互图 G_u 和 G_i , 构建过程中涉及到同质交互图阈值 K 的选取. 本小节通过实验说明预测误差和训练时间(每轮迭代耗时)随阈值 K 的变化情况, 并给出一种在兼顾效果和效率情况下选取最佳 K 值的方法. K 值对 GCN4RS 效果和效率的影响如图 4 所示.

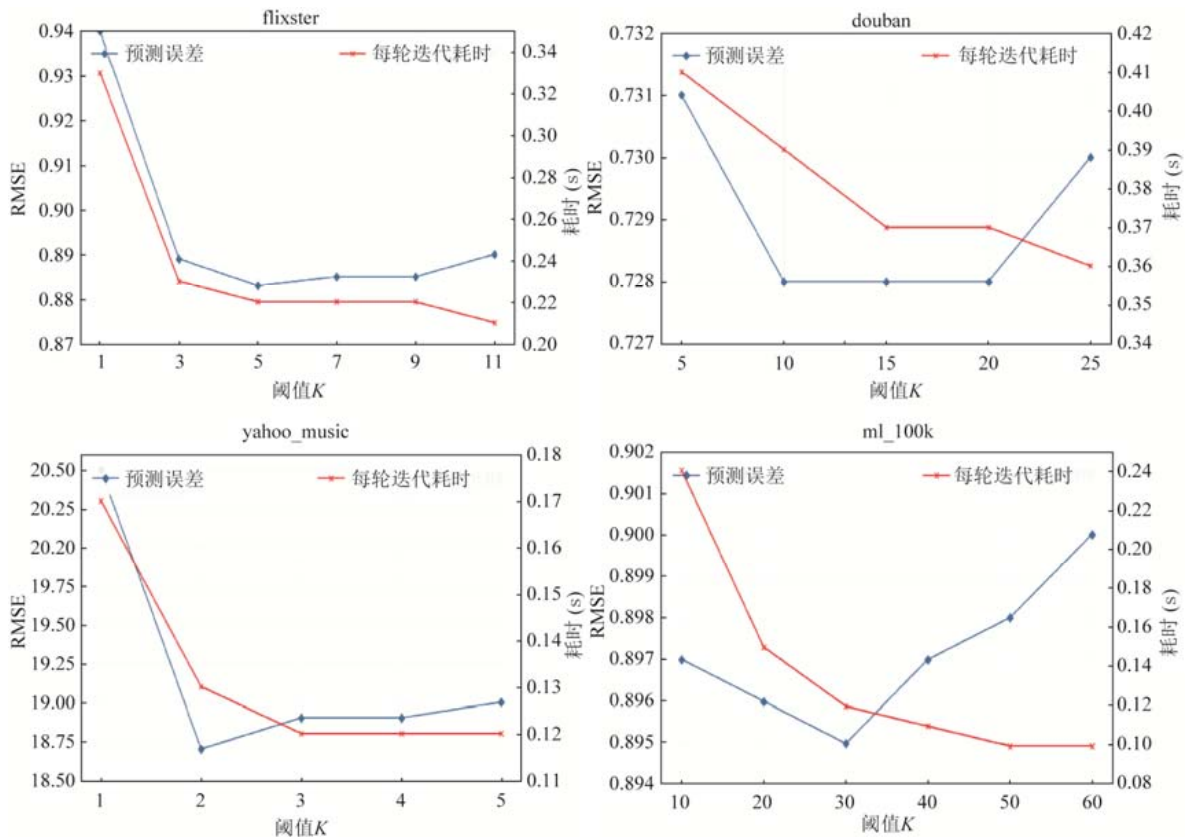


Fig.4 Selection of threshold K

图 4 阈值 K 的选取

由图 4 可知, 随着 K 的增加, 预测误差先降后增, 训练时间不断缩短. 这是因为, 若 K 过小, 一些不置信的链接被加入到同质交互图中, 一方面会引入噪声使效果下降, 另一方面会大幅度增加计算量而影响训练速度. 若 K 过大, 则同质交互图中的链接会变少, 训练加快, 但同质交互信息不足, 算法效果会有所下降.

预测误差和训练时间在 K 增加的过程中均存在斜率突变的拐点. 在拐点左侧, 预测误差和训练时间都快速降低; 在拐点右侧, 预测误差开始增加, 训练时间继续缩短, 但速度远低于拐点左侧. 可依据拐点位置选择阈值 K : 在拐点左侧, 预测误差过大, 不应选取这样的 K ; 在拐点右侧, 可权衡效果和效率后选择最合适的 K , 也可直接选取

拐点横坐标值作为 K . 本文为 Flixster、Douban、YahooMusic 和 MovieLens 数据集分别选取 5、15、2、30 作为阈值 K .

4.2.2 结果与分析

实验结果汇总见表 3. 所有对比方法均采用对应文献中的默认参数设置. GCNCF 的结果取自文献[14], 其未在 MovieLens 数据集上进行实验, 故未汇报结果.

其中, MC、GMC 和 GRALS 属于矩阵补全类算法, 其余算法均属于 GCN 类算法. 在 GCN 类算法中, GC-MC 和 GCN4RS-hetero 属于异质交互 GCN 类算法, RGCNN、GCNCF 和 GCN4RS-homo 属于同质交互 GCN 类算法. GCN4RS 是本文算法, 统一利用异质与同质信息. 我们将对比各类算法效果, 并给出相应理论分析.

Table 3 Experimental results

表 3 实验结果

	Flixster	Douban	YahooMusic	MovieLens
MC ^[24]	1.428	0.902	44.2	0.973
GMC ^[25]	1.411	0.878	40.4	0.996
GRALS ^[26]	1.245	0.833	38.0	0.945
GC-MC ^[17]	0.917	0.734	20.5	0.905
RGCNN ^[19]	0.926	0.801	22.4	0.929
GCNCF ^[14]	0.903	0.729	19.0	—
GCN4RS-Hetero	0.890	0.730	19.2	0.897
GCN4RS-Homo	0.941	0.814	19.0	0.953
GCN4RS	0.883	0.728	18.7	0.895

我们的算法 GCN4RS 在 4 个数据集上都取得了最优结果. 此外, 由实验结果可得到如下结论.

(1) 矩阵补全类算法 MC、GMC、GRALS 的效果显著差于 GCN 类算法. 矩阵补全类算法将用户与商品的交互视为矩阵, 重点挖掘线性相关与低秩信息; GCN 类算法重点挖掘交互图中的信息. 图相比于矩阵可刻画更多信息, 例如链接刻画相邻顶点间的联系, 拉普拉斯矩阵刻画图所有顶点间的整体联系, 链接密度刻画图中社区结构. 正是由于图具有矩阵无法比拟的强大表示能力, GCN 类算法的效果显著优于矩阵补全类算法, 这印证了第 3.4 节中的结论.

(2) 异质交互类算法和同质交互类算法的效果是可比的. 两类算法采用不同的图卷积方式, 但两者都将推荐系统建成图信号, 区别仅在于图信号的各成分指代的信息不同. 虽然同质交互类算法将评分信息视为矩阵而不是图, 但评分信息会以顶点信号的形式参与同质图上的图卷积运算, 使此类方法依然能够受益于图的强大表示能力.

(3) 除 YahooMusic 外, GCN4RS-Hetero 在其他数据集上的效果均显著优于 GCN4RS-Homo. 这是因为 GCN4RS-Homo 完全未利用评分信息, 而评分信息恰为推荐系统的核心. GCN4RS-Hetero 能够很好地利用评分信息, 但忽略了顶点相似性, 相比于 GCN4RS-Hetero, GCN4RS 的提升部分就来自对顶点相似性信息的利用. Yahoo Music 数据集的特点是评分密度很低(见表 2), 即异质交互信息很少, GCN4RS-Hetero 无充足的异质交互信息可用, 反而丢失了同质交互信息, 故 GCN4RS-Hetero 在此数据集上的效果略差于 GCN4RS-Homo.

(4) GCN4RS 的效果优于异质交互类和同质交互类算法. 相比于异质交互类算法, GCN4RS 利用了顶点间的相似性信息, 使相似顶点具有相近的嵌入向量表示, 在进行推荐时, 相似顶点更可能产生相近的行为, 这符合推荐系统基本假设; 相比于同质交互类算法, GCN4RS 用图而不是用矩阵来刻画评分信息, 能够更好地利用交互图谱域中蕴含的深层次连接信息, 不再局限于观测到的链接.

5 总结与展望

本文解决的问题是如何为推荐系统设计更合理的图卷积网络算法. 首先根据信息利用方式的不同, 将现有基于图卷积网络的推荐算法分类为异质顶点交互算法和同质顶点交互算法, 而两类方法都忽略了两类间的互助. 正是为了两者能够互惠互利, 本文提出了一种联合利用异质和同质交互图的图卷积网络算法. 真实数据集上的实验结果表明, 本文方法具有比现有方法更优的性能.

未来的改进方向是更合理地在图卷积操作中利用边权信息展开.GCN4RS 为每一级评分构建一个交互图,独立地利用这些交互图,最后将各个交互图上的信息融合并利用.然而,评分间存在有序关系,即 $5 > 4 > 3 > 2 > 1$,不同评分的交互图实际上是存在关联的,忽略此类关系会造成推荐性能的下降.因此,未来的工作应当探索如何在GCN4RS 中嵌入评分有序信息,可考虑借助有序回归(ordinal regression)等模型进行改进.

此外,在实际生产环境中,如何对图卷积操作进行改进,使其能够适应大规模图也是值得探索的问题.GCN4RS 中的图卷积采用基于矩阵乘法的谱域卷积,虽然能够有效挖掘图谱域中的信息,但在内存受限时难以进行运算.可考虑的改进方向是选用基于采样-聚合的卷积操作,进行分批次的、分布式的运算.

References:

- [1] LeCun Y, Bottou L, Bengio Y, *et al.* Gradient-based learning applied to document recognition. *Proc. of the IEEE*, 1998,86(11): 2278–2324.
- [2] LeCun Y, Kavukcuoglu K, Farabet C. Convolutional networks and applications in vision. In: *Proc. of the 2010 IEEE Int'l Symp. on Circuits and Systems*. IEEE, 2010. 253–256.
- [3] Dos Santos C, Gatti M. Deep convolutional neural networks for sentiment analysis of short text. In: *Proc. of the 25th Int'l Conf. on Computational Linguistics: Technical Papers*. 2014. 69–78.
- [4] Shuman DI, Narang SK, Frossard P, *et al.* The emerging field of signal processing on graphs: Extending high-dimensional data analysis to networks and other irregular domains. *IEEE Signal Processing Magazine*, 2013,30(3):83–98.
- [5] Bronstein MM, Bruna J, Lecun Y, *et al.* Geometric deep learning: Going beyond Euclidean data. *IEEE Signal Processing Magazine*, 2017,34(4):18–42.
- [6] Bruna J, Zaremba W, Szlam A, *et al.* Spectral networks and locally connected networks on graphs. In: *Proc. of the Int'l Conf. on Learning Representations*. 2014.
- [7] Defferrard M, Bresson X, Vandergheynst P. Convolutional neural networks on graphs with fast localized spectral filtering. In: *Advances in Neural Information Processing Systems*. 2016. 3844–3852.
- [8] Kipf TN, Welling M. Semi-supervised classification with graph convolutional networks. In: *Proc. of the Int'l Conf. on Learning Representations*. 2017.
- [9] Henaff M, Bruna J, LeCun Y. Deep convolutional networks on graph-structured data. *arXiv Preprint arXiv:1506.05163*, 2015.
- [10] Niepert M, Ahmed M, Kutzkov K. Learning convolutional neural networks for graphs. In: *Proc. of the Int'l Conf. on Machine Learning*. 2016. 2014–2023.
- [11] Chen J, Ma T, Xiao C. FastGCN: Fast learning with graph convolutional networks via importance sampling. In: *Proc. of the Int'l Conf. on Learning Representations*. 2018.
- [12] Hechtlinger Y, Chakravarti P, Qin J. A generalization of convolutional neural networks to graph-structured data. *arXiv Preprint arXiv:1704.08165*, 2017.
- [13] Ying R, He R, Chen K, *et al.* Graph convolutional neural networks for Web-scale recommender systems. In: *Proc. of the 24th ACM SIGKDD Int'l Conf. on Knowledge Discovery & Data Mining*. ACM, 2018. 974–983.
- [14] Jiang Y. Information fusion recommendation based on convolutional graph and neural collaborative filtering [MS. Thesis]. Changchun: Jilin University, 2018 (in Chinese with English abstract).
- [15] Qu Q, Yu HT, Huang RY. Spammer detection technology of social network based on graph convolution network. *Chinese Journal of Network and Information Security*, 2018,4(5):39–46 (in Chinese with English abstract).
- [16] Cai XD, Wang M, Liang XX, Chen Y. Community detection method based on graph convolutional network via importance sampling. *Journal of Zhejiang University (Engineering Science)*, 2019,(3):1–6 (in Chinese with English abstract).
- [17] Berg R, Kipf TN, Welling M. Graph convolutional matrix completion. *arXiv Preprint arXiv:1706.02263*, 2017.
- [18] Zheng L, Lu C T, Jiang F, *et al.* Spectral collaborative filtering. In: *Proc. of the 12th ACM Conf. on Recommender Systems*. ACM, 2018. 311–319.
- [19] Monti F, Bronstein M, Bresson X. Geometric matrix completion with recurrent multi-graph neural networks. In: *Advances in Neural Information Processing Systems*. 2017. 3697–3707.

- [20] Wu Y, Liu H, Yang Y. Graph convolutional matrix completion for bipartite edge prediction. In: Proc. of the Int'l Joint Conf. on Knowledge Discovery, Knowledge Engineering and Knowledge Management (KDIR). 2018. 51–60.
- [21] Hammond DK, Vandergheynst P, Gribonval R. Wavelets on graphs via spectral graph theory. Applied and Computational Harmonic Analysis, 2011,30(2):129–150.
- [22] Hinton GE, Salakhutdinov RR. Reducing the dimensionality of data with neural networks. Science, 2006,313(5786):504–507.
- [23] Kingma DP, Ba J. ADAM: A method for stochastic optimization. arXiv Preprint arXiv:1412.6980, 2014.
- [24] Candès EJ, Recht B. Exact matrix completion via convex optimization. Foundations of Computational Mathematics, 2009,9(6):717.
- [25] Kalofolias V, Bresson X, Bronstein M, *et al.* Matrix completion on graphs. arXiv Preprint arXiv:1408.1717, 2014.
- [26] Rao N, Yu HF, Ravikumar PK, *et al.* Collaborative filtering with graph information: Consistency and scalable methods. In: Advances in Neural Information Processing Systems. 2015. 2107–2115.

附中文参考文献:

- [14] 江原. 基于图卷积与神经协同过滤的融合信息推荐模型[硕士学位论文]. 长春: 吉林大学, 2018.
- [15] 曲强, 于洪涛, 黄瑞阳. 基于图卷积网络的社交网络 Spammer 检测技术. 网络与信息安全学报, 2018, 4(5): 39–46.
- [16] 蔡晓东, 王萌, 梁晓曦, 陈昀. 基于重要性抽样的图卷积社团发现方法. 浙江大学学报(工学版), 2019, (3): 1–6.



葛尧(1996—), 男, 山东泰安人, 硕士, CCF 学生会员, 主要研究领域为模式识别, 机器学习.



陈松灿(1962—), 男, 博士, 教授, 博士生导师, CCF 高级会员, 主要研究领域为模式识别, 神经计算.