

基于图的多源数据融合框架研究



匡广生^{1,2} 郭岩² 俞晓明² 刘悦² 程学旗²

1 中国科学院大学 北京 100049

2 中国科学院计算技术研究所 中国科学院网络数据科学与技术重点实验室 北京 100190

(kuangguangsheng@ict.ac.cn)

摘要 在给定的任务中分析各种数据时,目前大多数研究只针对单源数据进行分析,缺乏应用于多源数据的方法。但如今数据日益丰富,因此提出一种多源数据融合框架,用于融合多种网络平台数据。同一平台数据中包含文本与各种属性,同时不同平台的数据在内容与形式方面也存在很大差异。然而现有的网络信息挖掘方法大多仅使用同一平台中的部分数据进行分析,忽略了不同平台的数据之间存在的相互作用。因此文中提出一种数据融合框架,一方面,能基于图的强大表示能力融合同一平台不同类型的特征,从而提升单个平台的任务性能;另一方面能够利用不同平台的数据特征,使其相互补充,从而提升多个平台的任务性能。文中讨论的融合数据类型包括文本、时间、作者信息,这些特征涉及连续特征、离散特征以及非结构化特征。所提框架在事件分类任务上提升了F1值,验证了提出的多源数据框架的有效性。

关键词: 融合表示;多源数据;图融合

中图法分类号 TP391

Study on Multi-source Data Fusion Framework Based on Graph

KUANG Guang-sheng^{1,2}, GUO Yan², YU Xiao-ming², LIU Yue² and CHENG Xue-qi²

1 University of Chinese Academy of Sciences, Beijing 100049, China

2 Key Laboratory of Network Data Science & Technology, Institute of Computing Technology, Chinese Academy of Sciences, Beijing 100190, China

Abstract When analyzing various data in a given task, most of current researches only analyze single-source data and lack methods applied to multi-source data. But now data are becoming more abundant, therefore, this paper proposes a multi-source data fusion framework for fusing data from multiple network platforms. The data of the same platform contains text and various attributes, and there are also great differences in content and form among data of different platforms. Most existing network information mining methods only use part of the data in the same platform for analysis, and even ignore the interaction between the data of different platforms. Therefore, this paper proposes a data fusion framework, which can not only use more features of the same platform to improve the performance of a single platform, but also fuse the data features of different platforms to complement each other, thereby improving the performance of multiple platforms. This paper uses the task of event classification, and the abundant features effectively improve the F1 value, which verifies the effectiveness of the proposed multi-source data framework.

Keywords Fusion representation, Multi-source, Graph fusion

随着互联网的普及,海量的互联网内容纷涌而来。这些内容来自各种不同的网络平台,而且形式多种多样。以微博与新闻为例,两者的数据不仅在文本内容的风格上存在差异,而且由于微博属于新媒体,在信息的传播规律上两者也存在巨大差异。因此我们亟需更通用的数据表示方法来帮助我们充分利用多种平台的数据,从而更高效地完成各种分析任务。

首先,传统的数据表示主要集中在一种形式上,如文本形式的 doc2vec^[1]、网络结构形式的 Deepwalk^[2]等。但是即使

在同一个平台,其内部数据的构成也存在很大差异,比如微博的数据平台内既包含了文本信息、社交关系,又包含了用户的个人信息。我们将内容形式没有差异的平台定义成一种通道,例如把所有的新闻网站定义为新闻通道。目前的网络数据表示与分类研究较多基于单通道,并且大多仅利用单通道中的部分信息,所以多源信息融合是亟待解决的问题^[3]。如何融合同一通道内尽可能多的信息来进行后续的分析是本文要解决的第一个关键问题。

收稿日期:2020-11-02 返修日期:2021-03-18 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:国家重点研发计划(2017YFB0803302)

This work was supported by the National Key Research and Development Program of China(2017YFB0803302).

通信作者:郭岩(guoy@ict.ac.cn)

其次,对于不同通道,比如新闻与微博,它们虽然在各方面都存在差异,但是却共同构成了分析任务的对象,如网络事件等。如何对不同通道的数据进行融合对齐,利用不同通道的信息之间的相互促进作用来提高任务的准确率是本文要解决的第二个关键问题。

因此,本文的研究目标是融合多种数据源,提出基于图的融合框架,为后续的分析任务提供低维稠密数据表示。需要解决的问题主要有两个:1)需要融合同一通道内的多种数据特征,例如时间属性与文本特征;2)这些特征可能完全不一样,比如是否离散或者方差均值都不一样,从而导致融合难度较大。建立基于图的表示模型是由于图的表达能力非常强,节点与边都可以被赋予不同的涵义,即使差异非常大的特征,也可以通过多张图进行表示。因此我们将不同的特征都转化为图进行表示。

综上,本文的研究工作主要包括基于图来进行不同类型数据的表示、对同一通道数据进行融合以及将不同通道的数据进行对齐。最后我们利用寻找到的共同的表示方法来进行事件的识别。

本文的贡献主要有以下3点:1)提出了基于图的多通道数据融合模型;2)提出单通道与多通道的数据融合方法;3)用事件分类来验证框架的有效性。

1 相关工作

1.1 数据的图表示

数据表示算法包括无监督算法与监督算法。无监督算法将原始网络的信息表示转换为统一的低维语义空间的特征表示,同质的网络表示方法主要有谱方法、神经网络方法和矩阵分解方法;而监督算法主要是将神经网络算法引入各项任务。

基于谱方法的表示学习通过求解特定矩阵特征值和特征向量来得到节点的低维表示,线性空间的表示方法有线性判别分析,而流形空间的表示方法有非线性嵌入、非线性局部嵌入。基于矩阵分解的表示学习方法的主要思想是,利用给定的关系矩阵进行降维,并得到其低维表示。代表方法为GraRep^[4]。

对于上述两种方法,基于矩阵分解和谱方法的网络表示学习方法涉及特征向量计算、奇异值分解计算等复杂度较高的计算,并不适合大规模网络的实际应用。基于神经网络的表示学习方法利用神经网络的非线性特征提取能力,能够较好地挖掘网络节点间的复杂非线性关系。它利用神经网络建模节点间的网络结构信息,进而得到节点表示向量。Deepwalk^[2]方法首次将基于神经概率语言模型的词向量表示方法引入网络表示学习中。在Deepwalk算法的基础上,Line^[5]算法在节点采样时联合深度优先游走与宽度优先游走;Node2vec^[6]则利用有偏采样进行随机游走,使其能够表示出结构的相似性;而struct2vec^[7]引入了结构的相似性。异质网络考虑到边的有向性与节点的不同类型,以metapath2vec^[8]为代表的算法将随机游走变成随着有向边的路径游走,同时对节点概率归一化,考虑同种类型的节点,该异质网络还引入了注意力机制。

另一方面,监督算法可以依赖标注标签进行学习,因此通常比无监督算法的效果好。与自然语言和图像处理任务类似,图的学习中也引入了图卷积网络^[9-11]与图注意机制^[3,12],主要用于在图的邻居节点共享信息,极大地提高了图的学习能力。

1.2 数据的图融合

数据融合是有效提升相关任务准确率的手段。对于同一通道的数据融合算法,根据融合内容可以分为融合节点文本、融合节点属性、融合节点符号3种^[13]。

节点与节点文本的融合,可以认为是节点与节点中非结构化数据的融合。以TADW^[14]为代表,在矩阵形式的Deepwalk算法中,以目标矩阵分解形式引入文本内容,从而融合更多源信息。融合节点属性以SNE^[15]为代表,融合特征是节点本身包含的属性,比如节点与其他节点的连接数量。融合节点符号以MF^[16]方法为代表,在连边符号中,边可以代表用户之间的信任与否、朋友、敌人等关系,这个方法证明了关系中的弱结构平衡可以通过低秩矩阵分解得到。

除了按内容分类,还可以将通道融合分为前向融合与后向融合。前向方法在特征层面融合,后向方法在决策层融合。前向融合又叫早融合,与后向融合相比,能更好地捕捉不同通道的信息之间的相互作用。典型的通道融合就是跨媒体融合,大多采用前向融合方法,其目的是找到不同通道共同的语义空间。例如,coral^[17]算法将不同通道的隐层单元的方差进行对齐,从而寻找不同通道的共同语义空间,并将该空间用于相关的分析任务;对抗网络^[18]同样对不同来源的数据进行对抗学习。数据融合在城市研究^[19]、材料研究^[20]等各种领域有着广泛应用。

2 基于图的多源数据融合框架

本文提出的多源数据融合框架分成4个阶段,如图1的系统框架图所示。第一阶段是数据转化图的过程,不同的属性都可以转化成一张图,利用图来构建学习特征,即节点与边的信息,以便充分利用图的表达能力。第二阶段是单通道内的特征学习阶段,主要是将节点进行低维表示,同时还需要进行低维表示的聚合。第三阶段是对齐不同通道数据,即将通道内的表示向其他通道投影,得到统一的表示空间。第四阶段是分析任务目标的输出。

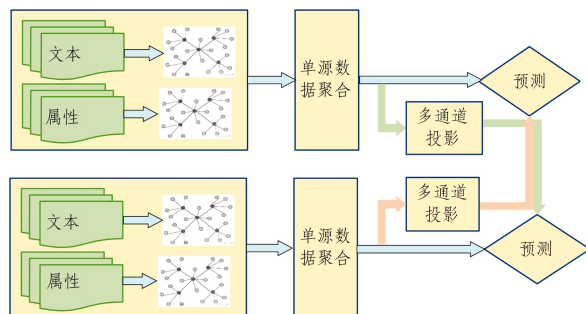


图1 系统框架图

Fig. 1 System frame diagram

2.1 图的构建

本文主要融合各个网络平台的数据,网络数据中含有文

本、发布者、发布时间等内容,需要对各个数据属性定义相似性。本文将数据分为3种类型,即连续结构化数据、离散属性数据以及非结构化数据。为不失一般性,本文选择文本、时间相关的属性、用户信息作为3种类型的代表,进行图的构建。图的质量至关重要。本文分别为以这3种类型为代表的图建立数据之间的相似性,从而建立高质量的图。

对于图的构建,属于同一类的聚集系数必然比属于不同类的聚集系数大,将不同属性或者特征转化成一张图,为后面的分类或者聚类任务做铺垫。

2.1.1 基于文本的图构建

对于文本这种非结构性的数据,节点之间的连接强度表示数据之间的相关性。在文本中,不同数据之间的连边可以基于文本的相似性来构建。主要流程如下:

(1)进行词预处理,包括分词、去停止词、去低频词等。

(2)建立文档-词矩阵,矩阵的每一项为逆文档频率。文档-词矩阵用 Jaccard 相似度计算文档相似度,从而得到文档-文档的矩阵,每一项即为边的权重。Jaccard 相似度的定义如下:

$$\text{sim}(D_i, D_j) = \frac{W_i \cap W_j}{W_i \cup W_j} \quad (1)$$

其中, D_i 是文档 i 的标记符号, W_i 为文档 i 的集合。由于文本可能涉及非常多的计算量,可以考虑使用分布式算法来计算边的权重。

2.1.2 基于时间属性的图构建

时间能很好地度量不同数据之间的相似性,时间越接近的内容,属于同一类的概率越大。因此需要度量时间远近用以表示数据之间的相似性,从而构建图。本文借鉴文献[21],用热力核函数对时间属性建图,有如下公式:

$$\text{sim}(D_i, D_j) = \exp\left(-\frac{\Delta t^2}{\sigma^2}\right) \quad (2)$$

其中, σ 为文档集合的时间方差, Δt 为两个文档之间的时间差,时间可以用天或者秒表示。本文对相似度做截断处理,若文档之间相似度低于阈值,则认为文档节点没有连接。经实验,阈值一般在 0.3~0.8 之间。

2.1.3 基于用户信息的图构建

本文也对用户信息进行建图,如果两条数据属于同一用户,则两条数据就有权值为 1 的连边。这种表示符合人们对数据的一般认识,因此离散数据的建立可以遵循这样的标准。此方法有可能导致图不连通,但是由于我们只对周围节点进行聚合,所以影响可以忽略不计。

得到不同特征的不同图之后,每个节点属性将同时被赋予所有的特征数据,这样即使某些数据不进行建图,也可以作为节点特征参与计算。

2.2 基于图的多通道融合算法

2.2.1 基于图注意机制的单通道融合

获得图及其节点属性之后,进入第二阶段,即节点表示,其通过融合并学习每个节点的特征得到,以便为后续任务服务。我们用图注意机制学习节点的特征。文献[1,3]的图注意过程是:对相邻节点经过全连接网络输出拼接,对其归一化

得到每个相邻节点的权重,然后更新本节点的隐层。

然而,每一个相邻节点在与本节点进行全连接输出的拼接过程中,参数太多,搜索范围太大,因此在重要性度量方面,本文直接用降维后的隐层单元相似性的向量积来计算节点相似度:

$$e_{ij} = W \vec{h}_i^T * W \vec{h}_j \quad (3)$$

该计算机制通过相邻节点的内积来衡量节点间的相似性。同时本文借鉴文献[3]的注意力机制,对计算得到的节点相似度进行归一化:

$$\alpha_{ij} = \text{softmax}_j(e_{ij}) = \frac{\exp(e_{ij})}{\sum_{k \in N_i} \exp(e_{ik})} \quad (4)$$

最终,得到最后的图网络输出:

$$\vec{h}_i' = \sigma\left(\sum_{j \in N_i} \alpha_{ij} W \vec{h}_j\right) \quad (5)$$

本文对建立的所有图使用多头注意机制,对每个图网络输出进行连接,从而得到第一层图网络的输出 \vec{h}^1 。第二层融合不同的图,本文的聚合信息方式为对每一个网络节点进行参数共享,从而得到融合输出:

$$\vec{h}_i^2 = W \vec{h}_i^1 \quad (6)$$

本文用事件分类对单通道的融合模型进行验证。事件分类本质上是一个多分类问题,因此单通道的损失函数为:

$$l_{\text{single}} = -\sum_k y_k \log\left(\frac{\exp(h_i^2)}{\sum_j \exp(h_j^2)}\right)_k \quad (7)$$

2.2.2 多通道融合表示

本文认为同一个标签下不同通道有相似的表达。因此,可以利用这种多通道的相似性来融合更多特征,从而利用通道之间的互补性来提高分类准确性。

这里以第一个通道分析其共同的语义空间。为了建模通道与通道之间的节点相互作用,将其他通道投影到第一通道空间,即其他通道的每一个节点输出投影到第一通道,因此可以对同一个标签空间进行参数共享。得到多通道的投影机制如下:

$$h_i^3 = h_i^2 + \sum_c (H * W_i)_c \quad (8)$$

其中, $H * W_i$ 为其他通道的数据输出对本通道的影响。对所有的输出标签进行累加即可得到多通道中此节点的综合输出。其他通道的投影机制类似。

输出包括原来单通道的输出和映射通道的输出,映射通道输出作为代价函数的正则项,其输出对最终的单通道有促进作用,此损失记为 l_{fusion} 。本文同样用事件分类来对多通道数据融合进行验证。因此本文总的代价函数为:

$$\text{cost} = l_{\text{single}} + \alpha * l_{\text{fusion}} \quad (9)$$

其中, α 为可调参数,经实验选择,这里取 0.1。

2.3 任务学习

本文的输出分为两部分,第一部分是单通道的训练输出,第二部分是投影部分,作为第一部分的约束。本文的数据融合所建模的是不同通道之间的相同目标标签下的共同语义空间,因此不是针对特定的具体任务,框架具有较好的泛化性。含多源在内的特征数据对任务的准确性有提升价值。值得注意的是,本文是在图结构的基础上进行分析,标签具有较好的

传播性,因此在少量标签下所提算法也能达到较好的性能。图融合算法的步骤如下。

(1)图的构建

1)对需要建图的属性,根据选择的相似度算法或者准则进行连边;

2)选择所有特征作为节点特征,并进行归一化。

(2)通道融合

1)根据 2.2.1 节进行单通道的融合;

2)根据 2.2.2 节进行多通道的融合。

(3)目标输出

输出最后的预测目标。

3 实验分析

3.1 实验数据及实验设置

本文的实验数据来自文献[22],具体来自 Flickr、YouTube 以及不同的新闻网站。该数据包含 23 874 个 Flickr 图像信息、10 678 条新闻报道信息以及 1 337 条 YouTube 视频信息,其中有不同的关键字,YouTube 的样例数据如表 1 所列。

表 1 样例数据

Table 1 Samples

| Title | Time | label |
|-----------------------------|------------|-------|
| Flood in Malaysia, Dec 2014 | 2015-01-10 | 0 |

本文抽取了数据集中前面的 15 个类标签,具体包括 6 659 个 Flickr 图像、2 606 个新闻报道、294 个 YouTube 内容。事件包括 2013—2015 年的全球自然灾害、政治经济热点事件。

图的建立过程如下:

(1)时间以天为单位,按照 2.1.2 节的方式建图;

(2)将数据的标题、描述与关键词组合成文本,建立基于文本的图;

(3)新闻有作者属性,对作者属性进行建图;

(4)对特征做最大最小值归一化。节点属性特征包含文本选择的关键词和平台的点赞数、转发数等。

本文建立的带权图如表 2 所列。

表 2 构建的图网络

Table 2 Graph network constructed

| 通道 | Flickr | Youtube | 新闻 |
|---------|------------|---------|-----------|
| 节点数 | 6 559 | 294 | 2 606 |
| 文本图边数 | 9 186 565 | 10 092 | 1 387 908 |
| 时间图边数 | 16 741 960 | 35 380 | 3 592 234 |
| 作者信息图边数 | 1 852 001 | — | — |

为了验证本文提出的融合框架的有效性,本文将事件分类任务用于本框架的目标学习中。事件分类任务是监督式分类的一种,由于一个事件涉及方方面面,同时对事件的报道涉及各个平台,各个平台对事件的反应各不一样,因此事件表现的内容形式都有可能不一样,此任务能比较有效地验证算法的有效性。

本文基于两种算法设计对比实验。一种是对特征使用逻辑回归与支持向量机算法。另一种是在本文构建的文本图、

时间图和作者属性图上使用 Line 和 Deepwalk^[2]算法,通过将 Line 和 Deepwalk 用于不同的图实验,来衡量图的质量,并将其作为融合算法的基线算法。Line 和 Deepwalk 的输出通过逻辑回归算法训练预测。

3.2 评价方法

本文采用单类的准确率 P 、召回率 R 、 $F1$ 作为评价指标。其计算方式如下:

$$P = \frac{tp}{tp + fp} \quad (10)$$

$$R = \frac{tp}{tp + fn} \quad (11)$$

$$F1 = \frac{2 * P * R}{P + R} \quad (12)$$

其中, tp 指预测与真实标签都属于某一事件; fp 指真实标签属于某个事件,而预测结果不属于; fn 指真实标签与预测都不属于某一类。

对实验数据含有多个标签评价来说,有 marco-F1 和 micro-F1 两个评价标准。macro-F1 先对所有的类求准确率与召回率的平均值,然后得到 $F1$ 值。micro-F1 则先计算每个类的 $F1$ 值,然后对各个类的 $F1$ 值取平均值。macro-F1 相比 micro-F1,对数据较少的类敏感,micro-F1 则对数据较多的类敏感,因此这两个指标能综合表明算法的整体性能。

3.3 实验结果及分析

本文的实验环境设置:linux 服务器 linux, GPU Tesla k80, python3 环境。

对单通道在含有 10% 的训练标签的情况下进行实验,对比结果如表 3 所列。

从 Line 和 Deepwalk^[2]算法与逻辑回归、SVM 算法两种非建图算法的对比来看,本文在建图策略与建图质量上表现很好, $F1$ 值有所提高,说明图的表达能力非常强。

从单通道融合的表现来看,一方面,与基线算法相比,本文的单通道融合表示算法略好。特别地,与只使用时间图相比,单通道内作者属性图与文本图分别与时间图融合后,都能大幅度提高算法的效果。这是因为单纯用数据中的时间做分类误差非常大,必须依赖特征或者其他图来进行矫正。因此完全依赖时间图的基线算法效果不理想。另一方面,与单通道的特征只作为节点属性融合的时间图与文本图相比,单通道的图融合效果比单图的融合效果有小幅提升。

从多通道融合效果来看,多通道融合涉及数据不平衡问题,因此在损失函数中,各个通道的权重对多通道结果有影响。由于 Flickr 数据量大于其他数据量,本文设置 Flickr 损失是其他损失的两倍。实验表明新闻的结果略微下降,Flickr 和 YouTube 提升明显。多通道融合方法能促进不同通道信息的相互作用,尤其对于数据稀疏的 YouTube 数据源来说,其正是其他如 Flickr 数据的补充,使得 YouTube 的事件分类性能得到较大提升。但同时多通道融合中新闻的 $F1$ 值略微下降,这是由于融合过程中同时对多通道进行学习,自身丰富的信息对其他通道起了促进作用,但其他通道稀疏的信息带来了一定的噪声,从而使 $F1$ 值略微下降。

图 2 给出了 $F1$ 值和训练节点标签比例的关系。从图 2 可以看出,即使标签很少,本文算法的泛化能力也非常强。

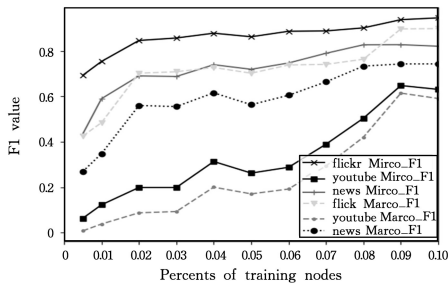


图2 F1值与训练节点比例的关系

Fig. 2 Relation between F1 value and trained nodes

综上所述,本文提出的多源数据融合框架,在单源数据上由于使用了图网络,高效地表征了数据之间的特征,并且可以有效融合单通道的信息;多维特征提高了分类能力,在多源数据方面,其能够通过多通道融合进一步提升通道间的任务学习能力。

表3 实验对比结果

Table 3 Experimental comparison results

| | | Flickr | YouTube | 新闻 |
|-----------------|----------|--------|---------|------|
| Logistic(只有特征) | macro-F1 | 0.22 | 0.14 | 0.14 |
| | micro-F1 | 0.65 | 0.41 | 0.50 |
| SVM(只有特征) | macro-F1 | 0.82 | 0.48 | 0.69 |
| | micro-F1 | 0.87 | 0.65 | 0.83 |
| Deepwalk(时间图) | macro-F1 | 0.53 | 0.18 | 0.21 |
| | micro-F1 | 0.78 | 0.43 | 0.65 |
| Deepwalk(文本图) | macro-F1 | 0.76 | 0.49 | 0.68 |
| | micro-F1 | 0.86 | 0.61 | 0.81 |
| Line(时间图) | macro-F1 | 0.34 | 0.29 | 0.33 |
| | micro-F1 | 0.85 | 0.76 | 0.80 |
| Line(文本图) | macro-F1 | 0.67 | 0.43 | 0.67 |
| | micro-F1 | 0.90 | 0.81 | 0.91 |
| Deepwalk(作者信息图) | macro-F1 | 0.90 | — | — |
| | micro-F1 | 0.90 | — | — |
| 单通道融合(作者信息图+) | macro-F1 | 0.96 | — | — |
| | micro-F1 | 0.97 | — | — |
| 单通道融合(时间图) | macro-F1 | 0.78 | 0.37 | 0.73 |
| | micro-F1 | 0.90 | 0.42 | 0.80 |
| 单通道融合(文本图) | macro-F1 | 0.76 | 0.50 | 0.78 |
| | micro-F1 | 0.87 | 0.57 | 0.85 |
| 单通道融合(时间图+文本图) | macro-F1 | 0.85 | 0.50 | 0.79 |
| | micro-F1 | 0.92 | 0.57 | 0.86 |
| 多通道融合(时间图+文本图) | macro-F1 | 0.90 | 0.59 | 0.75 |
| | micro-F1 | 0.95 | 0.63 | 0.82 |

结束语 本文提出了一种多源数据融合框架,该框架能够有效地融合单通道与多通道数据,提升网络分析任务的学习能力。该框架先建立数据的各个属性的图,用图表示数据,使得数据的表达与扩展性都非常强。该框架通过图的注意力机制学习单个特征下的低维表示,然后融合单通道的图表示,并用多通道空间相互投影对单通道结果进行对齐矫正。本文用事件分类验证了所提框架的有效性。

未来,我们需要设计更合理的通道之间的投影机制,以期得到更高质量的数据表示,进一步提升网络分析任务的性能。同时本文只在单机上进行了相关的实验,需要更多地考虑分布式学习的构建。

参考文献

[1] LE Q, MIKOLOV T. Distributed representations of sentences

and documents[C]//Proceedings of International Conference on Machine Learning. China: PMLR, 2014: 1188-1196.

- [2] 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: Association for Computing Machinery, 2014: 701-710.
- [3] ZHANG J, WANG Y, LI K H, et al. Multi-source Sensor Body Area Network Data Fusion Model Based on Manifold Learning [J]. Computer Science, 2020, 47(8): 323-328.
- [4] CAO S S, LU W, XU Q K. GraRep: Learning Graph Representations with Global Structural Information[C] // Proceedings of the 24th ACM International Conference on Information and Knowledge Management (CIKM '15). New York: Association for Computing Machinery, 2015: 891-900.
- [5] TANG J, QU M, WANG M, et al. LINE: Large-scale information network embedding [C] // International Conference on World Wide Web. 2015: 1067-1077.
- [6] GROVER A, LESKOVEC J. Node2vec: Scalable Feature Learning for Networks[C] // ACM Sigkdd International Conference on Knowledge Discovery & Data Mining. ACM, 2016.
- [7] RIBEIRO L F R, SAVERESE P H P, FIGUEIREDO D R. Struc2vec: Learning node representations from structural identity[C] // Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining. 2017: 385-394.
- [8] DONG M Y, CHAWLA N V, SWAMI A. metapath2vec: Scalable Representation Learning for Heterogeneous Networks[C] // the 23rd ACM SIGKDD International Conference. ACM, 2017.
- [9] DEFFERRARD M, BRESSON X, VANDERGHEYNST P. Convolutional neural networks on graphs with fast localized spectral filtering[J]. Advances in Neural Information Processing Systems, 2016, 29: 3844-3852.
- [10] KIPF T N, WELING M. Semi-supervised classification with graph convolutional networks[C] // 5th International Conference on Learning Representations. 2017.
- [11] SCARSELLI F, GORI M, TSOI A C, et al. The graph neural network model[J]. IEEE Transactions on Neural Networks, 2008, 20(1): 61-80.
- [12] VELIČKOVIĆ P, CUCURULL G, CASANOVA A, et al. Graph Attention Networks[C] // International Conference on Learning Representations. 2018.
- [13] LIU Z M. Research on network representation Learning Method based on heterogeneous information fusion [D]. Information Engineering University, Strategic Support Forces, 2018.
- [14] YANG C, LIU Z Y, ZHAO D L, et al. Network representation learning with rich text information[C] // Proceedings of the 24th International Conference on Artificial Intelligence (IJCAI'15). Buenos Aires, Argentina: AAAI Press, 2015 S: 2111-2117.
- [15] LIAO L, HE X, ZHANG H, et al. Attributed social network embedding[J]. IEEE Transactions on Knowledge and Data Engineering, 2018, 30(12): 2257-2270.
- [16] HSIEH C J, CHIANG K Y, DHILLON I S. Low rank modeling

- of signed networks[C]//Proceedings of the 18th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, Beijing, China, 2012:507-515.
- [17] SUN B, SAENKO K. Deep CORAL: Correlation Alignment for Deep Domain Adaptation [C] // Lecture Notes in Computer Science (ECCV 2016). Workshops, ECCV. Cham: Springer, 2016:443-450.
- [18] YU C, WANG J, CHEN Y, et al. Transfer learning with dynamic adversarial adaptation network [C] // 2019 IEEE International Conference on Data Mining (ICDM). IEEE, 2019: 778-786.
- [19] LIU J, LI T, XIE P, et al. Urban big data fusion based on deep learning: An overview [J]. Information Fusion, 2020, 53: 123-133.
- [20] ZHOU J, HONG X, JIN P. Information Fusion for Multi-Source Material Data: Progress and Challenges [J]. Applied Sciences, 2019, 9(17):3473.
- [21] YANG Z, LI Q, LU Z, et al. Dual structure constrained multi-modal feature coding for social event detection from flickr data [J]. ACM Transactions on Internet Technology (TOIT), 2017, 17(2):1-20.
- [22] LIN Z, YANG Z, SITU R, et al. Improving Maximum Classifier Discrepancy by Considering Joint Distribution for Domain Adaptation [C] // WISE 2018. Cham: Springer, 2018:253-268.



KUANG Guang-sheng, born in 1995, postgraduate. His main research interests include natural language processing and data fusion.



GUO Yan, born in 1974, Ph.D, associate researcher. Her main research interests include network information acquisition and so on.