

GRAND+: 可扩展的图形随机神经网络

冯文正¹, 董玉晓¹, 黄庭林¹, 尹子奇², 徐成¹³ Evgeny Kharlamov⁴, 唐杰¹

¹清华大学, ²北京理工大学

³腾讯公司, ⁴博世人工智能中心

fzw17@mails.tsinghua.edu.cn, {yuxiaod, tl091, jietang}@tsinghua.edu.cn, ziqiyin18@gmail.com

alexcheng@tencent.com, evgeny.kharlamov@de.boosch.com

ABSTRACT

图神经网络 (GNNs) 已被广泛用于图的半监督学习。最近的一项研究表明, 图随机神经网络 (GRAND) 模型可以为这个问题产生最先进的性能。然而, GRAND 很难处理大规模的图, 因为它的有效性依赖于计算上昂贵的数据增强程序。在这项工作中, 我们提出了一个可扩展和高性能的 GNN 框架--GRAND+, 用于半监督的图学习。为了解决上述问题, 我们在 GRAND+ 中开发了一个通用前推 (GFPush) 算法, 预先计算一个通用的传播矩阵, 并采用它以小批量的方式进行图数据增强。我们表明, GFPush 的低时间和空间复杂度使 GRAND+ 能够有效地扩展到大型图。此外, 我们在 GRAND+ 的模型优化中引入了一个置信度感知的一致性损失, 促进了 GRAND+ 的泛化优势。我们在七个不同规模的公共数据集上进行了广泛的实验。结果表明, GRAND+: 1) 能够扩展到大型图, 并且比现有的可扩展 GNN 花费更少的运行时间; 2) 在所有数据集上比全批和可扩展 GNN 都能提供一致的准确性改进。

CCS概念

• 计算方法学→半监督学习设置; -信息系统→社会网络。

关键字

图谱神经网络; 可扩展性; 半监督学习

ACM参考格式:

冯文正, 董玉晓, 黄庭林, 尹子奇, 徐成, Evgeny Kharlamov, 唐杰。2022. GRAND+: 可扩展的图随机神经网络-作品。在 2022 年 ACM 网络会议 (WWW'22) 论文集中, 2022 年 4 月 25-29 日, 虚拟活动, 法国里昂。ACM, New York, NY, USA, 11 pages. <https://doi.org/10.1145/3485447.3512044>

该代码可在 <https://github.com/wzfaha/GRAND-plus>。唐杰是通讯作者。

允许为个人或课堂使用本作品的全部或部分内容制作数字或硬拷贝, 不收取任何费用, 但拷贝不得以营利或商业利益为目的而制作或分发, 且拷贝首页须注明本通知和完整的引文。必须尊重 ACM 以外的其他人拥有的本作品的版权。允许摘录并注明出处。以其他方式复制, 或重新发表, 张贴在服务器上或重新分发到名单上, 需要事先获得特别许可和/或付费。请从 permissions@acm.org 申请许可。WWW'22, 2022 年 4 月 25-29 日, 虚拟活动, 法国里昂。

© 2022 年计算机协会。ACM ISBN 978-1-4503-9096-5/22/04... \$15.00

<https://doi.org/10.1145/3485447.3512044>

1 简介

图结构在我们的物理世界和虚拟世界中都很常见, 如社会关系、化学键和信息扩散等。现实世界图数据固有的不完整性激发了人们对图上半监督学习问题的巨大兴趣[19, 33]。到目前为止, 图神经网络 (GNNs) 已经被许多人认为是解决这个问题的事实上的方法[1, 11, 19, 28, 29]。简而言之, GNNs 利用数据样本之间的图结构来促进模型预测, 使它们能够产生比传统的半监督学习方法更突出的性能改进[32]。

然而, 基于 GNN 的半监督学习解决方案仍然存在挑战。值得注意的是, GNN 的泛化通常并不构成其优势, 因为它们中的大多数只使用监督损失来学习参数[9, 19, 28, 29]。这种设置使得模型容易对有限的标记样本进行过度拟合, 从而降低了对未见过的样本的预测性能。为了克服这个问题, 图形随机神经网络 (GRAND) [12] 为 GNNs 设计了图数据增强和一致性正则化策略。这些设计使其在 Cora、Citeseer 和 Pubmed 的半监督学习中比现有的 GNN 有明显的性能提升。

具体来说, GRAND 开发了随机传播操作, 以产生有效的结构数据增强。然后, 对它进行训练, 包括对已标记节点的监督损失和对未标记节点的不同扩增的一致性正则化损失。为了实现良好的图增强, GRAND 中的随机传播提出使用混序邻接矩阵来传播特征矩阵。传播基本上需要在每个训练步骤中对邻接矩阵进行幂级迭代, 这使得 GRAND 在扩展到大规模图时面临计算上的挑战。

为了解决这个问题, 我们提出了 GRAND+ 框架, 用于大规模的图上半监督学习。GRAND+ 是一种可扩展的 GNN 一致性正则化方法。在 GRAND+ 中, 我们引入了有效的近似技术, 以小批量的方式进行随机传播, 解决了 GRAND 的可扩展性限制问题。此外, 我们通过采用置信度感知损失来调节不同图数据增强之间的一致性, 从而改进 GRAND。这种设计稳定了训练过程, 并为 GRAND+ 提供了良好的泛化能力。具体来说, GRAND+ 包括以下技术:

• **广义的特征传播:** 我们提出一个广义的混序矩阵来进行随机特征传播。这种矩阵

它提供了一套可调整的权重来控制不同等级邻域的重要性,从而为处理复杂的现实世界图提供了一个灵活的机制。

- **高效的近似**: 受最近基于矩阵近似的GNNs[5, 8]的启发, GRAND+采用了一种近似方法——广义前推(GFPush)—有效地计算出

广义的传播矩阵。这使得GRAND+能够按

以小批量的方式形成随机传播和模型学习,为该模型提供了显著的可扩展性。

- **信心感知的损失**: 我们为GRAND+正则化框架设计了一个信心感知损失。这有助于在一致性训练中通过忽略高度不确定的未标记样本来过滤掉潜在的噪音,从而提高GRAND+的泛化性能。

我们在七个不同类型和规模的公共图数据集上进行了综合实验,以证明GRAND+的性能。总的来说,与10个GNN基线相比,GRAND+在三个基准数据集上产生了最好的分类结果,并在其他四个相对较大的数据集上超过了5个有代表性的可扩展GNN,而且效率很高。例如,GRAND+在Pubmed上达到了最先进的85.0%的准确性。在有1240万个节点的MAG-Scholar-C上,GRAND+比FastGCN和GraphSAINT快10倍,并且在运行时间相当的情况下,比PPRGo—以前在这个数据集上最快的方法—提供了4.9%的精度提升。

2 半监督的图形学习

2.1 问题

在基于图的半监督学习中,数据样本量或一个。

数据样本, $E \subseteq V \times V$ 是一组边,表示每对节点之间的关系-船舶。我们使用 $A \in \{0, 1\}^{|V| \times |V|}$ 来代表 G 的邻接矩阵,每个元素 $A(s, \square) = 1$ 表示 s 和 \square 之间存在一条边,否则 $A(s, \square) = 0$ 。

D 是对角线度矩阵,其中 $D(\square, \square) = \sum_{s \in V} A(s, \square)$ 用来表示添加了自环连接的图 \square 。相应的邻接矩阵为 $\tilde{A} = A + I$, 程度矩阵为

是 $\tilde{D} = D + I$

在这项工作中,我们关注的是分类问题,每个样本 s 都与1) 特征向量 $x_s \in \mathbb{R}^{|V| \times d}$ 和2) 标签向量 $y_s \in \mathbb{Y} \subseteq \{0, 1\}^{|V| \times C}$ 代表了类的数量。 $|V| \times C$, C 代表类的数量。在半监督设置中,只有有限的节点 $L \subseteq V$ 有观察到的标签 ($0 < |L| \ll |V|$), 其余节点的标签 $U = V - L$ 是未见的。半监督图学习的目的是根据图结构 \square 、节点特征 x 和观察到的标签 y_L ¹, 推断出未被标记的节点 U 的缺失标签 y_U 。

2.2 相关工作

图形神经网络 (GNNs) 已被广泛采用于半监督图形学习问题的修正。在这一部分,我们回顾了GNNs的进展,重点是它们对半监督图学习的大规模解决方案。

¹ 对于矩阵 $M \in \mathbb{R}^d \times \mathbb{R}$, 我们用 $M_{\square} \in \mathbb{R}^d$ 来表示其 \square -th 行向量, 并让 $M(\square, \square)$ 代表第1行和第1列的元素。

图卷积网络。图卷积网络 (GCN) [19] 将卷积操作概括为图。具体来说,第 \square 个GCN层定义为:

$$H^{(\square+1)} = \sigma(A^{(\square)} H^{(\square)} W^{(\square)}), \quad (1)$$

其中, $H^{(\square)}$ 表示第 \square 个隐藏节点的表示。

含有 H 的层 $(0) = x, \hat{A} = D^{-1/2} A D^{-1/2}$ 是对称的归一化

\square 的邻接矩阵, $W^{(\square)}$ 表示第三层的权重矩阵, $\sigma(\cdot)$ 表示激活函数。在实践中,这个图形卷积过程将被重复多次,最终的表征通常被送入逻辑回归层进行分类。

简化图卷积。通过仔细观察方程1,我们可以发现图卷积包括两个操作: **特征传播** $A^{(\square)} H^{(\square)}$ 和 **非线性变换** $\sigma(\cdot)$ 。Wu等人[29]通过去除隐藏层中的非线性变换来简化这一过程。由此产生的简化图卷积 (SGC) 被表述为:

$$Y^* = \text{softmax}(A^{(\hat{N})} x w), \quad (2)$$

其中, $A^{(\hat{N})} x$ 被认为是一个简化的 N 层图对 x 的卷积, w 是指用于分类的可学习权重矩阵, Y^* 表示模型的预测。

具有混合顺序传播的GNNs。正如Li等人[23]所指出的,根据马尔科夫链收敛定理, $A^{(\hat{N})} x$ 将随着 N 的增加而收敛到一个固定点,即过度平滑问题。为了解决这个问题,一种典型的方法[5, 20, 21] 建议使用一个更复杂的混序矩阵来处理特征传播。例如, APPNP[20] 采用了截断的每——归纳的网页排名 (ppr) 矩阵 $\Pi^{\text{ppr}} = \square \alpha (1 - \square)^n A^n$, 其中

的超参数

征兆 $\square = 0$

$\alpha \in (0, 1)$ 表示远距离传输概率,使模型能够保留

本地信息,即使在

$N \rightarrow +\infty$ 。

可扩展的GNNs。大体上,有三类方法被提出来用于使GNNs可扩展: 1) 节点采样方法

采样策略,以加快速度归特征的速度。

聚合程序。代表性的方法包括 Graph-SAGE[14]、FastGCN[7] 和 LADIES[34]; 2) 图的分割方法试图将原始大图分割成几个小的子图,并在子图上运行GNNs。这类方法包括 Cluster-GCN[10] 和 GraphSAINT[31]; 3) 矩阵近似方法遵循SGC[29]的设计,将特征传播和非线性转换解耦,并利用一些近似方法来加速特征传播。所提出的GRAND+框架与基于矩阵逼近的方法高度相关,如PPRGo[5]和GBP[8]。我们将在第3.6节中分析它们的区别。

3 大+框架

在本节中,我们简要回顾了图随机神经网络 (GRAND), 并介绍了其用于大规模半监督图学习的可扩展解决方案GRAND+。

3.1 图形随机神经网络

最近, Feng等人[12]介绍了用于半监督节点分类的图神经网络 (GRAND)。GRAND是一个GNN

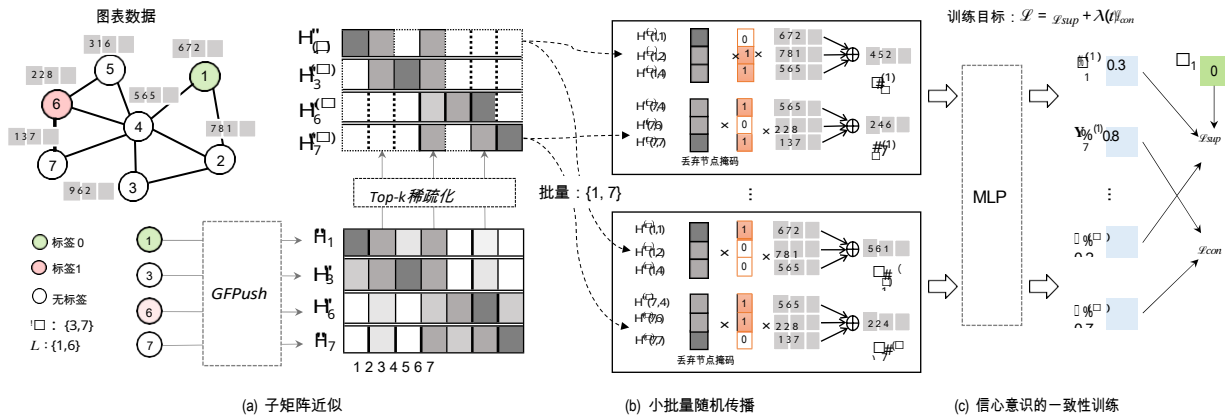


图1: GRAND+的说明。 (a) GRAND+采用Generalized Forward Push (GFPush)和Top-k稀疏化, 对 L_{UU} 中的节点进行传播矩阵 Π 的相应行近似。(b) 然后利用得到的稀疏化行近似进行小批量随机传播, 为批量中的节点生成增强值。(c) 最后, 将计算出的特征增量送入MLP进行置信度感知的一致性训练, 该训练同时采用监督损失 \mathcal{L}_{sup} 和置信度感知的一致性损失 \mathcal{L}_{con} 进行模型优化。

一致性正则化框架, 优化了不同增量中未标记节点的预测一致性。

具体来说, 它设计了随机传播--一种混合顺序的传播策略--来实现图数据的增强。首先, 用DropNode--一种dropout的变体, 随机地丢弃节点特征 x 。然后, 由此产生的被破坏的特征矩阵用混序矩阵在图上传播。而不是PPR矩阵。

进行传播。形式上, 随机传播策略被模拟为: $\bar{x} = \Pi \bar{z} - \text{diag}(z) \cdot x$, $z_i \sim \text{Bernoulli}(1 - \delta)$ 。GRAND使用一个平均集合矩阵 $\bar{x} = \frac{1}{M} \sum_{m=1}^M \bar{x}^{(m)}$ 。

where $z \in \{0, 1\}^{|U|}$ 表示从Bernoulli($1 - \delta$)中抽取的随机DropNode掩码, 而 δ 代表DropNode概率。这样一来, 每个节点丢失的信息就由其邻域来补偿。在图形数据的同质性假设下, 产生的矩阵 x 可以被看作是原始特征矩阵 x 的有效数据增强。由于DropNode的随机性, 这种方法在理论上可以为每个节点产生指数级的许多增强值。

在GRAND的每个训练步骤中, 随机传播程序被执行了 M 次, 导致 M 增强的特征矩阵 $\{\bar{x}^{(m)} | 1 \leq m \leq M\}$ 。然后将所有增强的特征矩阵输入MLP, 得到 M 的预测结果。在优化过程中, GRAND的训练既包括对已标记数据的标准分类损失, 也包括对未标记节点集 U 的额外一致性正则化损失[4], 即: 、

$$\frac{1}{M} \sum_{m=1}^M \frac{1}{|U|} \sum_{s \in U} \|\hat{y}_s^{(m)} - \hat{y}_s\|^2, \quad \hat{y}_s = \text{MLP}(\bar{x}^{(m)}), \quad \bar{x}^{(m)} = \Pi \bar{z}^{(m)} - \text{diag}(z) \cdot x, \quad z_i \sim \text{Bernoulli}(1 - \delta) \quad (4)$$

其中 $\hat{y}_s^{(m)}$ 是MLP对节点 s 的预测概率。当使用 $\bar{x}^{(m)}$ 作为输入。一致性损失提供了一个额外的regu-

通过强制神经网络对不同的无标签数据进行类似的预测, 从而达到拉长的效果。通过运行穹顶传播和一致性正则化, GRAND实现了比传统的GNN有更好的泛化能力[12]。

GRAND 的可扩展性。 在实践中, 当 n 较大时, 对增强矩阵 $A^{(n)}$ 的第 n 幂进行计算是不可行的[25]。为了避免这个问题, GRAND采用幂级迭代法直接计算整个增强特征矩阵 x (在公式3中), 即迭代计算并求出 $\hat{A}^{(n)}$ 和 $A^{(n)} - \text{diag}(z) \cdot x$ 的乘积为 $0 \leq n < N$ 。这个过程是用稀疏密集矩阵乘法实现的, 具有线性时间

复杂度。在每个训练步骤中产生不同的特征增强。因此, T 训练步骤的总复杂度变成了 $O(N \cdot M \cdot (|V| + |E|))$, 这在处理大型图时是非常昂贵的。

3.2 GRAND+的概述

我们提出GRAND+来实现基于图的半监督学习的可扩展性和准确性。它遵循GRAND的一般一致性正则化原则, 并包括一些技术, 使其可以扩展到大型图, 同时保持GRAND的灵活性和概括能力。

简而言之, 我们在GRAND+中开发了一种高效的近似算法--广义forward push (GFPush), 以预先计算传播矩阵所需的行向量, 并以小批量的方式进行随机传播, 而不是用功率迭代来传播特征。该程序的时间复杂度由一个预定义的超参数控制, 避免了GRAND所面临的可扩展性限制。此外, GRAND+采用了一种新的置信度感知损失来实现一致性正则化, 这使得表现比GRAND好。

传播矩阵。 在GRAND+中, 我们提出以下基因

侵蚀的混序矩阵用于特征传播:

$$\Pi = \frac{1}{N} W - P^n, \quad P = \text{diag}(p) \quad (5)$$

其中 $\Pi = \Pi_n$, $\Pi_n \geq 0$, P 是行规范化的邻接性矩阵。与GRAND中使用的传播矩阵不同的是

和其他GNN, Π 的形式采用了一组可调整的权重。通过调整 $\Pi \leq N$ 来融合不同顺序的邻接矩阵, 操纵其重要性。

的不同阶数的邻域, 以适应现实世界中不同结构属性的不同图形。

训练管道。为了实现快速训练, GRAND+放弃了直接计算整个增强特征矩阵 X 的功率迭代方法, 而是为每个节点单独计算每个增强特征向量。理想情况下, 节点 v 的增强特征向量 x_v 的计算方法是:

$$\bar{x}_v = \frac{1}{|\mathcal{N}(v)|} \sum_{v' \in \mathcal{N}(v)} z_{v'} - \Pi(s, v) - x_{v'}, z_{v'} \sim \text{Bernoulli}(1 - \Pi(s, v)). \quad (6)$$

这里我们用 Π_s 来表示 Π 的行向量, 对应于节点 s , Π_s 被用来表示 Π 的第 s 个元素。这种范式允许

我们在每个训练步骤中只为一批节点生成增强的特征, 从而使我们能够使用高效的小批量梯度下降法进行优化。

然而, 在实践中很难计算出 Π 的确切形式。技巧。为了解决这个问题, 我们开发了几种有效的方法

来逼近GRAND+中的 Π_s 。近似过程由两个阶段组成。在第一阶段, 我们提出了一种有效的方法——通用前推法(GFPush)来计算行向量的误差约束近似值。在第二阶段, 我们采用顶层稀疏化策略来截断元素, 使得到的稀疏化行近似于

$\Pi^{(k)}$ 被用来计算 x_v , 作为 Π_s 的替代(公式6)。

为了提高效率, 需要预先计算出相应的行

训练中使用的所有节点的近似值。除了标记的节点, GRAND+还要求无标签的节点执行一致的

在训练过程中, GRAND+采用了正则化技术。为了进一步提高效率, GRAND+没有使用全部的 U , 而是对一个较小的未标记节点的子集 $U' \subseteq U$ 进行一致性正则化。如图1所示, GRAND+的训练流水线包括

三个步骤:

- **子矩阵近似。**我们得到一个稀疏化的行近似, 通过GFPush对每个节点 $s \in L \cup U$ 进行建模 $\Pi^{(k)}$ 和顶部- k 稀疏化。由此产生的稀疏化子矩阵被用于支持随机传播。
- **小批量的随机传播。**在每个训练步骤中, 我们从 $L \cup U'$ 中抽出一批节点, 用近似的行向量为该批中的每个节点生成多个增强值。
- **信心意识到的一致性训练。**我们将增强的特征送入MLP以获得相应的预测, 并通过监督损失和置信度感知的一致性损失对模型进行优化。

3.3 子矩阵近似

广义前推(GFPush)。可以看出, 行归一化的邻接矩阵 A^{-1} 也是 Π 上的反向随机行走过渡概率矩阵 P^A , 其中行向量 p_v 表示从节点 v 开始的随机行走过渡概率。基于这一事实, 我们提出一种有效的算法

称为Generalized Forward Push (GFPush), 以近似于行向量。

的 $r_v = \Pi v - \Pi(s, v) - x_v$, 误差有限。GFPush的灵感来自于

前推[2]算法用于逼近个性化的PageRank向量, 同时具有更高的灵活性, 能够逼近广义的混序矩阵 Π 。GFPush的核心思想是模拟一个 N 步的随机行走概率扩散过程, 从 s 用一系列的剪枝操作进行acceleration。为了达到这个目的, 我们应该在每一步 n ($0 \leq n \leq N$) 保持一对向量: 1) 储备向量 $q^{(n)} \in \mathbb{R}^{|V|}$, 表示在第 n 步保留的概率质量; 2) 残余向量 $r^{(n)} \in \mathbb{R}^{|V|}$, 代表超过步骤 n 的扩散概率质量。

算法1显示了GFPush的伪代码。在开始时, $e_s^{(0)} = 1$, $e_v^{(0)} = 0$, 对于 $v \neq s$, $r_v^{(0)}$ 和 $q_v^{(0)}$ 都被初始化为指向向量 $e^{(0)}$, 其中 Π_s 的概率质量为1。其他储备和残留物向量(即 $r^{(n)}$ 和 $q^{(n)}$, $1 \leq n \leq N$)被设置为0。然后, 该算法以 Π 步迭代更新储备和残留值向量。

在第 n 迭代中, 算法对满足 $r_v^{(n-1)} > d_v - m_v$ 的节点 v 进行推送操作(算法1的第5-9行)。这里 $d_v = D(v, \Pi)$ 代表 v 的度数, m_v 是一个预定的阈值。在推送操作中, v 的残差 $r_v^{(n-1)}$ 是均匀的

分布到它的邻域, 并将结果存储到 v 的残余向量 $r_v^{(n)}$ 中。同时, 储备向量 $q^{(n)}$ 也被更新为与 $r^{(n)}$ 相同。在完成对 v 的推送操作后, 我们将 $r_v^{(n-1)}$ 重置为0, 以避免重复更新。

为了更直观地了解这一程序, 我们可以观察到

$r_v^{(n-1)} / d_v$ 是随机行走的条件概率。从 v 到邻接节点 u , 条件是它到达 v 时有上一步的概率 $r_v^{(n-1)}$ 。因此, 每次推送操作

在 Π 上, 可以被看作是一个单步随机行走的概率扩散

从 v 到其邻域的过程。为了保证效率, GFPush只对残差值大于 $d_v - m_v$ 的节点 v 进行推动操作。因此, 当第 n 迭代完成后, $q^{(n)}$ 可以被看作是对 Π 步数的随机行走的迹描述。

观念向量 p^v 。而 $\Pi - \Pi_s = \Pi - \Pi_s$ 相应地被认为是作为算法所返回的 Π_s 的近似值。

理论分析。我们有以下定理, 关于时间复杂度、内存复杂度和近似度的界限

GFPush的错误。

定理1。算法1的时间复杂度为 $O(N \cdot l_{max})$, 内存复杂度为 $O(N \cdot l_{max})$, 并返回 Π_s 作为 Π 的近似值, 其误差界限为: $\|\Pi - \Pi_s\|_1 \leq \frac{1}{2} (|E| + |\Pi|) \cdot l_{max}$ 。

证明。见附录A.2。

定理1表明, GFPush的近似精度和运行成本与 Π 呈负相关。在实践中, 我们可以使用 Π 来控制效率和近似精度之间的权衡。

顶层稀疏化(Top-k Sparsification)。为了进一步降低训练成本, 我们对 Π 进行顶部稀疏化。在这个过程中, 只保留了 Π 中最大的元素, 其他条目被设置为0。这样, 该模型只考虑了最多的

算法1: GFPush

输入 : 自环增强图 G , 传播步骤 N , 节点 \square , 阈值 m , 权重系数 $\square_n, 0 \leq n \leq N_0$.

输出: 节点 \square 的过渡向量 \square 的近似值。

```

1  $q^{(n)} \leftarrow 0$  对于  $n = 1, \dots, \square$ ;  $r^{(0)} \leftarrow e^{(\square)} (e^{(\square)} = 1, e^{(\square)} = 0$  对于  $v \neq \square$ )。
2  $q^{(n)} \leftarrow 0$  对于  $n = 1, \dots, \square$ ;  $q^{(0)} \leftarrow e^{(\square)}$ 。
3 对于  $\square = 1$ :  $\square$ 
do
4   如果存在节点 $v$ , 且  $r^{(n-1)} > d_{\square} \square_m \square$  做
5   对于每个  $u \in N_{\square}$  做
6   /*  $N_{\square}$  是  $\square$  在图  $\square$  中的邻居集。 */
7    $r_{\square}^{(n)} \leftarrow r_{\square}^{(n)} + r_{\square}^{(\square-1)} / d_{\square}$ 。
8    $q_{\square}^{(n)} \leftarrow r_{\square}^{(n)}$ 。
9   结束
10   $r_{\square}^{(\square-1)} \leftarrow 0$ 。
11 /* 对  $\square$  进行推操作。 */
12 结束
13 返回  $\square_{\square}$ 。

```

在随机传播中, 每个节点的重要邻域, 根据聚类假设[6], 这仍然被认为是有效的。类似的技术也被PPRGo[5]所采用。我们将在第4.5节中实证检验 \square 的效果。

并行化。在GRAND+中, 我们需要对 $L \cup U'$ 中的所有节点进行行向量近似。可以很容易地检查出, 在GFPush中, 不同的行近似是相互独立计算的。因此, 我们可以启动多个工作器, 同时对多个向量进行近似。这个过程是这样实现的

在我们的实现中采用多线程编程。

3.4 小批量随机传播

GRAND+采用 Π 的稀疏行近似值来执行

以小批量的方式进行随机传播。具体来说, 在

在第19个训练步骤中, 我们随机抽取一批有标签的节点 L 来自 \square , 而一批未标记的节点 U' 来自 U' 。然后, 我们计算节点 $s \in L \cup U'$ 的增强特征向量 X_{\square} :

$$X_{\square} = \sum_{\square \in N_{\square}} z_v \cdot \Pi^{(\square)}(s, \square) - X_v, z_v \sim \text{Bernoulli}(1 - \square), \quad (7)$$

其中 $N_{\square}^{(\square)}$ 表示 Π_{\square} 的非零指数, $X_v \in \mathbb{R}^{d_h}$ 是节点 v 的特征向量。在每个训练步骤中, 我们生成 \square

重复这个过程 M 次, 就可以得到增强的特征向量 $\{X_{\square}^{(m)} | 1 \leq m \leq \square\}$ 。让 $b = |\square| + |\square|$ 表示批次大小。那么每个批次的时间复杂度以 $O(\square \cdot \square \cdot \square)$ 为界。

可学习表征的随机传播。在公式7中, 增强的特征向量 X_s 是用原始特征 x 计算的。然而, 在一些实际应用中(如图像或文本分类), x 的维度可能非常大, 这将产生巨大的计算成本。为了缓解这个问题, 我们可以采用线性层将每个 x 转化为低维的隐藏表示 $H_v \in \mathbb{R}^{d_h}$ 。首先, 然后用 H 进行随机传播:

$$X_{\square} = \sum_{\square \in N_{\square}} z_v \cdot \Pi^{(\square)}(s, \square) - H_v, H_v = X_v - W^{(0)}, \quad (8)$$

其中 $W^{(0)} \in \mathbb{R}^{d_f \times d_h}$ 表示可学习的变换矩阵。这样, 该程序的计算复杂度就降低到 $O(\square \cdot \square \cdot \square)$, 其中 $d_h \ll d_f$ 表示 H 的维度。

预测。在训练期间, 增强的特征向量 $X^{(\square)}$

被送入一个MLP模型, 以获得相应的输出:

$$Y_{\square} = \text{MLP}(X_{\square}^{(m)}, \Theta), \quad (9)$$

其中 $Y_{\square}^{(m)} \in [0, 1]^C$ 表示 \square 的预测概率。 Θ 代表MLP的参数。

3.5 信心意识的一致性训练

GRAND+在训练过程中采用监督分类损失和一致性正则化损失来优化模型参数。监督下的损失被定义为以下方面的平均交叉熵

标记的节点的多个增量:

$$L_{sup} = - \frac{1}{|L|} \sum_{s \in L} \sum_{m=1}^M Y_s \cdot \log(Y_{\square}^{(m)}). \quad (10)$$

信心意识的一致性损失。受半监督学习的最新进展启发[4], GRAND采用了附加的一致性损失来优化未标记数据的多次增强的预测一致性, 这在提高泛化能力方面被证明是有效的。GRAND+也遵循了这一思想, 同时采用了新的置信度感知的一致性损失来进一步提高有效性。

具体来说, 对于节点 $s \in U'$, 我们首先通过取其 M 预测概率的平均值来计算分布中心, 即:

$\bar{Y}_{\square} = \sum_{m=1}^M Y_{\square}^{(m)} / M$ 。然后我们将 τ [17]的技巧应用于 \bar{Y}_{\square} 来“猜”出节点 \square 的伪标签 \tilde{Y}_{\square} 。从形式上看, 节点 j 的第 s 类上的猜测概率是通过以下方式获得的:

$$\tilde{Y}_{\square}(\square, \square) = Y_{\square}(\square, \square)^{\frac{1}{1+\tau}} \cdot \frac{1}{\sum_{\square=1}^C Y_{\square}(\square, \square)^{\frac{1}{1+\tau}}} \quad (11)$$

其中, $0 < \tau \leq 1$ 是一个超参数, 用于控制尖锐度。

猜测的伪标签。随着 \square 的减少, \tilde{Y}_{\square} 被强制变为

更加尖锐, 并最终收敛于单热分布。那么, 在无标签节点批上的置信度感知的一致性损失 L_{\square} 被定义为:

$$L_{\square} = \frac{1}{|U'|} \sum_{s \in U'} I(\max(\bar{Y}_s) \geq \gamma) \sum_{\square=1}^M D(Y_{\square}^{(m)}, \tilde{Y}_{\square}), \quad (12)$$

其中 $I(\max(\bar{Y}_{\square}) \geq \gamma)$ 是一个指示函数, 如果 $\max(\bar{Y}_{\square}) \geq \gamma$ 成立则输出1, 否则输出0。 $0 \leq \gamma < 1$ 是一个预定的阈值。 $D(p, \square)$ 是一个距离函数, 衡量 p 和 \square 之间的分布差异。这里我们主要考虑 D 的两个选项: \square_2 距离和KL分歧。与GRAND使用的一致性损失相比(参见公式4), L_{\square} 最大的优势在于它在优化过程中只考虑由阈值 τ 决定的“高度自信”的无标签节点。这种机制可以通过过滤掉不确定的伪标签来减少潜在的训练噪声, 进一步提高模型的实际性能。结合 L_{\square} 和 L_{sup} , 可得模型优化的最终损失被定义为:

$$L = L_{\square} + L_{sup}(\square), \quad (13)$$

算法2: GRAND+

```

输入 : 图  $G$ , 特征矩阵  $X \in \mathbb{R}^{V \times df}$ , 有标签的节点集  $L$ , 无标签
        的节点集  $U$  和观察标签  $y \in \mathbb{R}^{|U| \times \Omega}$ 。
输出: 分类概率  $\hat{Y} \in \mathbb{R}^{|U| \times \Omega}$ 。
1 从  $U$  中抽取一个无标签节点的子集  $U'$ 。
2 对于  $s \in L \cup U'$  做
3    $\Pi_s \leftarrow \text{GFPush}(G, s)$ 。
4   通过对  $\Pi_s$  进行顶- $\square$  稀疏化得到  $\pi^{(\square)}$ 。  $\sim$ 
   /* 用泛化的方式逼近行向量。 */
5 结束
6 for  $\square = 0 : \square_{\text{do}}$ 
7   对一批有标签的节点  $\square \in L$  和一批无标签的节点进行采样
    $U_i \subseteq U'$ 。
8   对于  $s \in L \cup U_i$  做
9     对于  $\square = 1 : \square_{\text{do}}$ 
10      生成增强的特征向量  $X_{\square}$ 。  $\sim$  (用方程7生成增强的
      特征向量  $X_{\square}$ 。)
11      用  $\hat{Y}^{(\square)} = \text{MLP}(X_{\square}^{(\square)}, \theta)$  预测类的分布。  $\sim$ 
12      结  $\hat{Y}^{(\square)}$ 
13      束
14      通过公式 10 计算  $L_{\text{sup}}$ , 通过公式 12 计算  $L \square \square_n$ 。通过小
      批量梯度下降来更新参数  $\theta$ :  $\theta = \theta - \nabla_{\theta} (L_{\text{sup}} + \lambda L_{\text{con}})$ 。
15      /* 用提前停止的方式停止训练。 */
       $\hat{Y}^{(\square)}$ 
16 结束
17 推断分类概率  $\hat{Y} = \text{MLP}(\Pi - (1 - \square) - X, \theta)$ 。
18 返回  $\hat{Y}^{(\square)}$ 。

```

其中, $\square(\square)$ 是一个线性预热函数[13], 随着训练步骤 \square 的增加, 它从 0 线性增加到最大值 $\square \square$ 。

模型推断。训练结束后, 我们需要推断出未标记节点的预测结果。在推理过程中, GRAND+ 采用幂级迭代的方式计算未标记节点的精确预测结果:

$$\hat{Y}^{(\square)}(\square) = \sum_{\square=0}^N n^w(\square) - (1 - \square) - X, \theta, \quad (14)$$

其中我们用 $(1 - \square)$ 重新标定 X , 以使其与训练中使用的 DropNode 扰动特征的期望值。

请注意, 与 GRAND 不同的是, 上述功率迭代过程只需要在 GRAND+ 中执行一次, 而计算

成本在实践中是可以接受的。与获得预测值相比

与训练中的 GFPush 相比, 这种推理策略在理论上可以提供更准确的预测。算法2显示了 GRAND+ 的整个训练和推理过程。

3.6 模型分析

复杂度分析。我们对 GRAND+ 不同学习阶段的时间复杂度进行了详细分析。根据定理1, 近似阶段 (算法2的第2-5行) 的复杂度为 $O((|U'| + |L|) \cdot N / \square \square_a \square)$ 。至于训练阶段 (算法2的第6-15行), T 训练步骤的总复杂度为 $O(k \cdot b \cdot M \cdot \square)$, 这对大图来说实际上是有效的, 因为 b 和 \square 通常远小于图的大小。推理阶段 (算法2第17行) 的复杂度为 $O((|V| + |E|) \cdot N)$, 与节点和边的数量之和成线性关系。

GRAND+ vs. PPRGo 和 GBP。与 GRAND+ 类似, PPRGo[5] 和 GBP[8] 也采用矩阵逼近的方法来扩展 GNNs。然而, GRAND+ 在几个关键方面与这两种方法不同。PPRGo 通过使用 Forward Push[2] 来扩大 APPNP 的规模。

表1: 数据集的统计数据。

数据集	节点	边缘	课堂	特点
Cora	2,708	5,429	7	1,433
Citeseer	3,327	4,732	6	3,703
公共医学杂志	19,717	44,338	3	500
琥珀-CS	593,486	6,217,004	18	100
睿迪特 (Reddit) 公司	232,965	11,606,919	41	602
Amazon2M	2,449,029	61,859,140	47	100
MAG Scholar-C	40,544,560	265,340,004	8	2,704,240

近似于 PPR 矩阵。与 PPRGo 相比, 由于采用了广义的传播矩阵 Π 和 GFPush 算法, GRAND+ 在实际应用中更加灵活。GBP 也拥有这个通过使用广义的 PageRank 矩阵[22] 进行特征分析, 来确定优点。

传播。然而, 它直接近似于传播的图, 通过双向传播[3] 对原始特征的结果进行处理, 其

计算复杂度与原始特征尺寸呈线性关系, 因此难以处理具有高维特征的数据集。此外, 与 PPRGo 和 GBP 的设计不同的是, 在一般的监督分类问题上, GRAND+ 会对其进行统计分析。

通过采用随机传播和一致性正则化来提高泛化能力, 对半监督设置进行了重大改进。

4 实验

4.1 实验设置

基线。在我们的实验中, 我们将 GRAND+ 与五个最先进的全批次 GNNs-GCN[19], GAT[28], APPNP[20], GCNII[9] 和 GRAND[12], 以及五个有代表性的可扩展 GNN--FastGCN[7], GraphSAINT[31], SGC[29], GBP[8], 和 PPRGo[5]。对于 GRAND+, 我们实现了三个变体, 分别是

传播矩阵 Π 的不同设置 (参考方程5):

- GRAND+ (P): 截断的 ppr 矩阵 $\Pi^{\text{PPR}} = \sum_{\square=0}^N \alpha (1 - \square)^{\square} \Pi^{\square}$ 。
- GRAND+ (A): 平均集合矩阵 $\Pi^{\text{平均}} = \frac{\sum_{\square=0}^{\square_{\text{do}}} \Pi^{\square}}{\square_{\text{do}} + 1}$ 。
- GRAND+ (S): 单阶矩阵 $\Pi^{\text{单}} = \Pi^{\square_{\text{do}}}$ 。

数据集。实验在七个不同规模的公共数据集上进行, 包括三个广泛采用的基准图--Cora、Citeseer 和 Pubmed[30], 以及四个相对较大的图--AMiner-CS[12]、Reddit[14]、Amazon2M[10] 和 MAG-Scholar-C[5]。对于 Cora、Citeseer 和 Pubmed, 我们使用公共数据拆分 [19, 28, 30]。对于 AMiner-CS、Reddit、Amazon2M 和 MAG-Scholar-C, 我们使用 $20 \times \# \text{classes}$ 节点进行训练, $30 \times \# \text{classes}$ 节点进行验证, 其余节点进行测试。表1中总结了相应的统计数据。更多关于设置和可重复性的细节可以在附录A.1中找到。

4.2 基准数据集的结果

为了评估 GRAND+ 的有效性, 我们在 Cora、Citeseer 和 Pubmed 上将其与 10 个 GNN 基线进行比较。按照社区惯例, 基线模型在这三个基准上的结果取自以前的工作[9, 12, 28]。对于 GRAND+, 我们用随机种子进行了 100 次试验, 并报告了平均精度和相应的标准差。

表2：基准的分类准确率（%）。

类别	方法	科拉	馨予	医学博士
全批次的GNNs	GCN	81.5 ± 0.6	71.3 ± 0.4	79.1 ± 0.4
	GAT	83.0 ± 0.7	72.5 ± 0.7	79.0 ± 0.3
	APNP	84.1 ± 0.3	71.6 ± 0.5	79.7 ± 0.3
	GCNII	85.5 ± 0.5	73.4 ± 0.6	80.3 ± 0.4
	盛大的快速GCN	85.4 ± 0.4	75.4 ± 0.4	82.7 ± 0.6
可扩展的GNNs	快速GCN	81.4 ± 0.5	68.8 ± 0.9	77.6 ± 0.5
	GraphSAINT	81.3 ± 0.4	70.5 ± 0.4	78.2 ± 0.8
	SGC	81.0 ± 0.1	71.8 ± 0.1	79.0 ± 0.1
	英镑	83.9 ± 0.7	72.9 ± 0.5	80.6 ± 0.4
	PPRGo	82.4 ± 0.2	71.3 ± 0.3	80.0 ± 0.4
我们的方法	GRAND+ (P)	85.8 ± 0.4	75.6 ± 0.4	84.5 ± 1.1
	Grand+ (a)	85.5 ± 0.4	75.5 ± 0.4	85.0 ± 0.6
	Grand+ (s)	85.0 ± 0.5	74.4 ± 0.5	84.2 ± 0.6

表3：大型图形的准确率（%）和运行时间（s）。

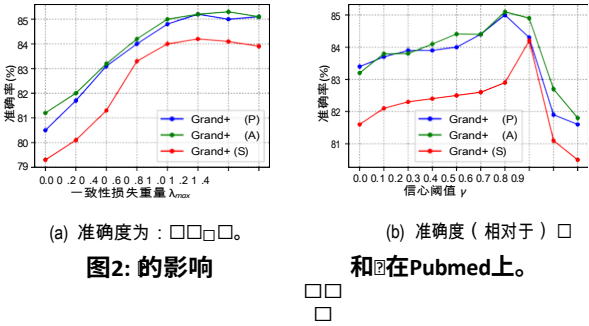
方法	AMiner-CS Acc 术 语	Reddit Acc 报 道	亚马逊2.5 亿美元 Acc 访 谈	MAG。 帐户 术 语
盛大的	53.1±1.1 750	OOM -	OOM -	OOM -
快速GCN	48.9±1.6 69	89.6±0.6 158	72.9±1.0 239	64.3±5.6 4220
GraphSAINT	51.8±1.3 39	92.1±0.5 39	75.9±1.3 189	75.0±1.7 6009
SGC	50.2±1.2 9	92.5±0.2 31	74.9±0.5 69	- >24h
英镑	52.7±1.7 21	88.7±1.1 370	70.1±0.9 280	- >24h
PPRGo	51.2±1.4 11	91.3±0.2 233	67.6±0.5 160	72.9±1.1 434
GRAND+ (P)	53.9±1.8 17	93.3±0.2 183	75.6±0.7 188	77.6±1.2 653
GRAND+ (A)	54.2±1.7 14	93.5±0.2 174	75.9±0.7 136	80.0±1.1 737
GRAND+ (S)	54.2±1.6 10	92.8±0.2 62	76.2±0.6 80	77.8±0.9 483

的试验。结果显示在表2中。我们可以看到，最好的GRAND+变体在三个数据集上的表现一直优于所有的基线。值得注意的是，GRAND+（A）在Pubmed上比GRAND提高了2.3%（绝对差异）。在Cora（85.8±0.4 vs. 85.4±0.4）和Citeseer（75.6±0.4 vs. 75.4±0.4）上，GRAND+（P）比GRAND的改进也具有统计学意义（经t检验，p值为0.01）。这些结果表明，GRAND+取得了强大的泛化性能。

4.3 大图上的结果

为了证明GRAND+的可扩展性，我们进一步将其与四个大图上的五个可扩展的GNN基线进行比较，即AMiner-CS、Reddit、Amazon2M和MAG-Scholar-C。请注意，MAG-Scholar-C的特征维度是巨大的（即每个节点2.8M的特征）。为了使GRAND+能够处理这个问题，在随机传播之前加入了一个可学习的线性层，将高维节点特征转化为低维的隐藏向量（参见公式8）。为了进行公平的比较，我们对所有方法进行了仔细的超参数选择（参见附录A.1）。我们对每个模型进行了10次随机拆分，并报告了它在试验中的平均精度和平均运行时间（包括预处理时间、训练时间和推理时间）。结果总结在表3中。

我们从两个角度来解释表3的结果。首先，结合表2的结果，我们注意到GRAND+的三个变体在这些数据集中表现出很大的差异：



在Cora和Citeseer上，GRAND+（P）比GRAND+（A）和GRAND+（S）取得了更好的结果；在Pubmed、Reddit和MAG-Scholar-C上，GRAND+（A）超过了其他两个变体；在AMiner-CS和Amazon2M上，GRAND+（S）得到了最佳的分类结果。这表明传播矩阵在这个任务中起着关键作用，并进一步表明GRAND+可以通过调整广义混序矩阵n来灵活地处理不同的图。

其次，我们观察到GRAND+在准确率上一直超过所有的基线方法，并在四个数据集上获得有效的运行时间。重要的是，在最大的图MAG-Scholar-C上，GRAND+可以在10分钟左右成功地完成训练和预测，而SGC和GBP需要24小时以上才能完成，因为这两种方法在预处理步骤中被设计为直接支持高维的原始特征。与FastGCN和GraphSAINT相比，GRAND+（S）分别实现了8倍和12倍的加速。与过去该数据集上最快的模型PPRGo相比，GRAND+ (S)在运行时间相当的情况下，准确率提高了4.9%。这些结果表明GRAND+在大图上的扩展性很好，并进一步强调了其卓越的性能。

我们还报告了GRAND在AMiner-CS上的准确性和运行时间。请注意，由于内存不足的错误，它不能在其他三个大数据集上执行。我们可以看到，GRAND+在AMiner-CS上的运行时间比GRAND快40倍以上，证明了所提出的近似技术在提高效率方面的有效性。

4.4 概括性的改进

在本节中，我们定量研究了所提出的置信度感知的一致性损失□□□□对模型的生成能力的好处。在GRAND+中，□□□□主要由两个超参数主导：置信度阈值γ（参照公式12）和最大一致性损失权重□□（参照公式13）。

我们首先分析了γ和□□□□对GRAND+的分类性能的影响。具体来说，我们调整了γ和□□□□的值。

λ_max，其他超参数固定，观察GRAND+在测试集上的准确性变化。图2说明了Pubmed数据集的结果。从图2（a）可以看出，当λ_max从0增加到0.8时，准确性明显提高。当□□□□大于0.8时，精度趋于稳定。这表明，一致性损失可以真正促进GRAND+的性能。从图2（b）中，我们观察到，当γ为0时，模型的性能从γ的扩大中获益。

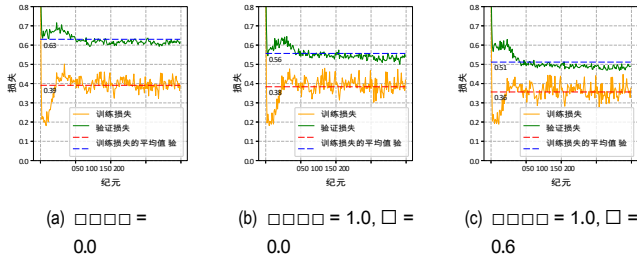


图3: Pubmed上的训练和验证损失。

小于0.7, 这突出了信心的重要性。

机制。如果 λ 设置得过大(即 >0.7), 性能就会下降, 因为在这种情况下, 太多的未标记样本被忽略, 削弱了一致性正则化的效果。

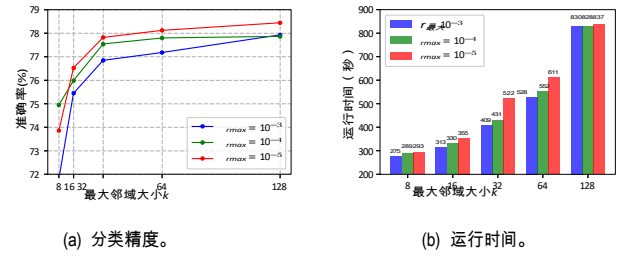
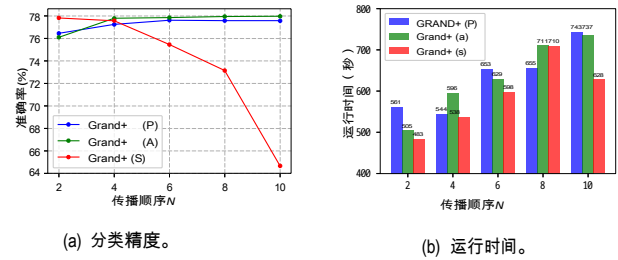
图2显示了置信度感知的一致性损失可以显著提高模型的性能。我们通过分析训练集和验证集的交叉熵损失, 进一步研究它对泛化能力的好处。这里我们用泛化差距[16]来衡量模型的泛化能力, 即训练损失和验证损失之间的差距。一个较小的泛化

化差距意味着该模型具有更好的泛化能力。图3报告了GRAND+ (A) 在Pubmed上的训练和验证损失。我们可以看到, 当我们在训练过程中不使用一致性损失($\lambda_{max}=0$)时, 泛化差距相当大, 表明存在严重的过拟合问题。而当我们把 λ 改为1.0时, 差距就会变小。当我们把 λ 和 γ 设置为适当的值(即 $\lambda_{max}=1.0$, $\gamma=0.6$), 泛化差距进一步减小。这些观察结果表明, 建议的一致性训练和信心机制确实有助于GRAND+的泛化能力。

4.5 参数分析

阈值 r_{ma} 和邻域尺寸 k 。 GRAND+使用GF-Push和top- k 稀疏化来逼近 n 的多个行向量, 以进行小批量的随机传播(参见第3.4节)。这个过程的近似误差主要受两个超参数的影响--GF-Push的阈值 r_{ma} 和稀疏化的最大邻域尺寸 k 。我们进行了详细的实验, 以更好地了解 k 和 r_{ma} 对模型的准确性和运行时间的影响。图4显示了GRAND+ (S) 在不同的 k 和 r_{ma} 对MAG-Scholar-S的相应结果。C. 我们可以看到, 当 r_{ma} 变小时, 精度和运行时间都会增加, 这与定理1的结论相吻合。而 k 则有相反的效果--随着 k 的增加, 精度和运行时间都会增加。有趣的是, 当 k 从128下降到32时, 运行时间减少了一半, 而准确度只下降了2%。这证明了顶层稀疏化策略的有效性, 它可以在几乎不牺牲精度的情况下实现大幅加速。

传播顺序 α 。 我们研究了在使用不同传播矩阵时, 传播顺序 α 对GRAND+的影响。图5显示了三种GRAND+变体在MAG-Scholar-C上的分类性能和运行时间, 其中 α 的数值不同。我们可以看到, 当 $\alpha=2$ 时, GRAND+ (S)取得了更好的成绩。

图4: GRAND+与 λ 和 γ 在MAG-Scholar-C上。图5: 传播顺序 α 对MAG-Scholar-C的影响。

比GRAND+(P)和GRAND+(A)的精度更高, 运行时间更快。

(A). 然而, 随着 α 的增加, 由于过度平滑问题, GRAND+ (S) 的精度急剧下降, 而GRAND+ (P)和GRAND+ (A)不存在这个问题, 并受益于较大的传播顺序。另一方面, 增加 N 会增加模型的运行时间。在实际应用中, 我们可以灵活地调整传播矩阵和 N 的值, 以取得理想的效率和效果。

5 结论

我们提出了GRAND+, 一个可扩展和高性能的GNN框架, 用于基于图的半监督学习。GRAND+的优势包括可扩展性和泛化能力, 而现有的最先进的解决方案通常只具有这两者中的一个。为此, 我们在实现泛化性能时遵循了GRAND的一致性调节原则, 同时对其进行了大幅扩展, 以实现可扩展性, 并保留甚至超越了GRAND的灵活性和泛化能力。为了实现这些, GRAND+利用广义混淆矩阵进行特征传播, 并使用我们的计算方法广义前推(GFPush)对其进行有效计算。此外, GRAND+还采用了一种新的置信度感知的一致性损失来实现更好的一致性训练。大量的实验表明, GRAND+不仅在基准数据集上获得了最好的表现, 而且在具有数百万节点的数据集上, 也取得了比现有可扩展GNN更高的性能和效率。在未来, 我们希望探索更精确的近似方法来加速GNNs。

鸣谢

该工作得到了国家自然科学基金委杰出青年学者(61825602)和清华-博世联合ML中心的支持。

参考文献

- [1] Sami Abu-El-Hajja, Bryan Perozzi, Amol Kapoor, Hrayr Harutyunyan, Nazanin Alipourfard, Kristina Lerman, Greg Ver Steeg, and Aram Galstyan. 2019. Mix-hop: 通过稀疏化邻域混合的高阶图卷积架构. *icml'19* (2019)。
- [2] Reid Andersen, Fan Chung, and Kevin Lang. 2006. 使用 pagerank 向量进行局部图分区. In *FOCS'06*. IEEE, 475-486.
- [3] Siddhartha Banerjee and Peter Lofgren. 2015. 马尔科夫模型中的快速双向概率估计. *NeurIPS* (2015)。
- [4] David Berthelot, Nicholas Carlini, J. Ian Goodfellow, Nicolas Papernot, Avital Oliver, and Colin Raffel. 2019. MixMatch: 一个半监督的整体方法学习. *NeurIPS* (2019), 5050-5060.
- [5] Aleksandar Bojchevski, Johannes Klicpera, Bryan Perozzi, Amol Kapoor, Martin Blais, Benedek Rózemberczki, Michal Lukasik, and Stephan Günnemann. 2020. 用近似的 pagerank 来扩展图神经网络. In *KDD'20*. 2464-2473.
- [6] Olivier Chapelle, Bernhard Schölkopf, and Alexander Zien. 2006. *半监督学习*. 麻省理工学院出版社. <https://doi.org/10.7551/mitpress/9780262033589.001>. 0001
- [7] 陈杰, 马腾飞, 和曹晓. 2018. FastGCN: Fast Learning with Graph Convolutional Networks via Importance Sampling. *ICLR* (2018)。
- [8] 陈明, 魏哲伟, 丁波林, 李亚良, 袁烨, 杜晓勇, 温继荣. 2020. 通过双向传播的可扩展图神经网络. *NeurIPS* (2020).
- [9] 陈明, 魏哲伟, 黄增丰, 丁波林, 和李亚良. 2020. 简单和深入的图卷积网络. 在 *ICML'20*. pmlr, 1725-1735.
- [10] 蒋伟林, 刘玄清, 思思, 李阳, Samy Bengio, 和谢祖瑞. 2019. Cluster-GCN: An Efficient Algorithm for Training Deep and Large Graph Convolutional Networks. In *KDD'19*.
- [11] Ming Ding, Jie Tang, and Jie Zhang. 2018. 图上的半监督学习与生成式对抗网. *cikm'18* (2018).
- [12] 冯文正, 张杰, 董玉晓, 韩瑜, 奕焕波, 徐倩, 杨强, 叶夫根尼-哈拉莫夫, 和唐杰. 2020. 图随机神经网络, 用于图上的半监督学习. *NeurIPS* (2020).
- [13] Priya Goyal, Piotr Dollár, Ross Girshick, Pieter Noordhuis, Lukasz Wesolowski, Aapo Kyrola, Andrew Tulloch, Yangqing Jia, and Kaiming He. 2017. 准确的、大型的 minibatch sgd: 在 1 小时内训练 imagenet. *arXiv 预印本 arXiv:1706.02677* (2017)。
- [14] Will Hamilton, Zhitao Ying, and Jure Leskovec. 2017. 归纳表征 大图上的学习. *NeurIPS* (2017), 1025-1035.
- [15] Sergey Ioffe and Christian Szegedy. 2015. 批量归一化: 通过减少内部协变量偏移来加速深度网络训练. *ICML* (2015)。
- [16] Shirish Nitish Keskar, Dheevatsa Mudigere, Jorge Nocedal, Mikhail Smelyanskiy, and Tak Peter Ping Tang. 2017. On Large-Batch Training for Deep Learning: Generalization Gap and Sharp Minima. *ICLR* (2017).
- [17] Sohn Kihyuk, Berthelot David, Li Chun-Liang, Zhang Zizhao, Carlini Nicholas, Ekin Cubuk D., Kurakin Alex, Zhang Han, and Raffel Colin. 2020. FixMatch: 简化具有一致性和信心的半监督学习. *NeurIPS* (2020).
- [18] P. Diederik Kingma 和 Lei Jimmy Ba. 2015. Adam: A Method for Stochastic Optimization. *ICLR* (2015)。
- [19] N. Thomas Kipf and Max Welling. 2017. 用图卷积网络进行半监督分类. *ICLR* (2017)。
- [20] Johannes Klicpera, Aleksandar Bojchevski, and Stephan Günnemann. 2019. Predict then Propagate: 图神经网络满足个性化的 PageRank. *ICLR* (2019).
- [21] Johannes Klicpera, Stefan Weissenberger, and Stephan Günnemann. 2019. Diffusion 改善了图的学习. *arXiv 预印本 arXiv:1911.05485* (2019)。
- [22] Pan Li, Eli Chien, and Olgica Milenkovic. 2019. 优化用于种子扩展社区检测的广义寻根方法. *arXiv 预印本 arXiv:1905.10881* (2019)。
- [23] Qimai Li, Zhichao Han, and Xiao-Ming Wu. 2018. 深入了解图卷积网络的半监督学习. 在 *AAAI'18*。
- [24] Razvan Pascanu, Tomas Mikolov, and Yoshua Bengio. 2013. 论训练递归神经网络的难度. *icml* (2013), 1310-1318.
- [25] Jiezhong Qiu, Yuxiao Dong, Hao Ma, Jian Li, Kuansan Wang, and Jie Tang. 2018. 网络嵌入作为矩阵分解: 统一 deepwalk、line、pte 和 node2vec. In *WSDM'18*. 459-467.
- [26] Arnab Sinha, Zhihong Shen, Yang Song, Hao Ma, Darrin Eide, Paul Bo-Jun Hsu, and Kuansan Wang. 2015. 微软学术服务 (MAS) 概述和应用. *WWW (配套卷)* (2015), 243-246.
- [27] 唐杰, 张静, 姚利民, 李娟子, 张丽, 和苏中. 2008. Arnet-miner: 学术社会网络的提取和挖掘. 在 *KDD'08* 中。
- [28] Petar Velickovic, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua Bengio. 2018. 图注意力网络. *ICLR* (2018)。
- [29] Felix Wu, Amauri Souza, Tianyi Zhang, Christopher Fifty, Tao Yu, and Kilian Weinberger. 2019. Simplifying graph convolutional networks. *ICML* (2019)。
- [30] 杨志林, William W Cohen, 和 Ruslan Salakhutdinov. 2016. 重新审视与图嵌入的半监督学习. *ICML* (2016)。
- [31] Hanqing Zeng, Hongkuan Zhou, Ajitesh Srivastava, Rajgopal Kannan, and Viktor Prasanna. 2020. GraphSAINT: 基于图形抽样的归纳学习方法. *ICLR* (2020).
- [32] Dengyong Zhou, Olivier Bousquet, Thomas N Lal, Jason Weston, and Bernhard Schölkopf. 2004. 具有局部和全局一致性的学习. *NeurIPS* (2004), 321-328.
- [33] Xiaojin Zhu, Zoubin Ghahramani, and John D Lafferty. 2003. 使用高斯场和谐波函数的半监督性学习. *ICML* (2003).
- [34] Difan Zou, Ziniu Hu, Yewen Wang, Song Jiang, Yizhou Sun, and Quanquan Gu. 2019. 用于训练深度和大图卷积网络的层依赖重要性采样. *NeurIPS* (2019)。

A 附录

A.1 实施说明

A.1.1 运行环境。实验在一台Linux服务器上运行, 该服务器配备英特尔(R)至强(R)CPU Gold 6420 @ 2.60GHz, 376G内存和10个NVIDIA GeForce RTX 3090TI-24GB。

Python的版本是3.8.5。

A.1.2 实施细节。我们用C++实现GFPush, 并使用OpenMP来执行并行化。我们使用Pytorch来实现GRAND+的训练过程, 并使用pybind⁵来创建近似模块的Python绑定。在GRAND+和其他基线中, 我们使用BatchNorm[15]和梯度剪裁[24]来稳定模型的训练, 并采用Adam[18]进行优化。

A.1.3 数据集详情。本论文共使用了7个数据集, 即Cora、Citeseer、Pubmed、AMiner-CS、Reddit、Amazon2M和MAG-Scholar-C。我们对Cora、Citeseer和Pubmed的预处理脚本是参照Pane-toid⁶的代码实现的。按照[19, 28, 30]中使用的实验设置, 我们在第4.2节中报告的Cora、Citeseer和Pubmed的结果中使用随机种子进行了100次试验。AMiner-CS是由Feng等人[12]基于AMiner引文网络[27]构建的。在AMiner-CS中, 每个节点代表一篇论文, 边是引文关系, 标签是论文的研究主题。Reddit是由Hamilton等人[14]发表的, 其中每个节点代表Reddit社区中的一个帖子, 一个图链代表两个帖子被同一个用户评论过。任务是预测每个帖子的类别。Amazon2M是由Chiang等人发表的[10], 其中每个节点是一个产品, 每个边表示两个产品一起购买, 标签代表产品的类别。MAG-Scholar-C由Bojchevski[5]基于微软学术图谱(MAG)[26]构建, 其中节点指的是论文, 边表示论文之间的引用关系, 特征是论文摘要的词包。

对于AMiner-CS、Reddit、Amazon2M和MAG-Scholar-C, 我们用20×#类用于训练, 30×#类节点用于验证, 其余节点用于测试。对于Aminer、Reddit和MAG-Scholar-C, 我们对每个类随机抽取相同数量的节点, 训练时每类20个节点, 验证时每类30个节点。对于Amazon2M, 我们从整个数据集中统一抽取所有的训练和验证节点, 因为有些类的节点数少于20。对于这些数据集, 我们报告了随机分割的10个追踪的平均结果。

A.1.4 超参数选择。对于表2的结果, 我们在验证集上调整了GRAND+的超参数, 并使用最佳配置进行预测, 其他基线方法的结果取自以前的工作[9, 12, 28]。对于表3-5的结果, 我们对GRAND+和其他GNN基线(即FastGCN、GraphSAINT、SGC、GBP和PPRGo)都进行了详细的超参数搜索。对于每一次搜索, 我们用随机种子进行3次实验, 并选择在验证集上获得最佳平均精度的超参数配置。然后, 我们用选定的配置训练模型。

⁵ <https://github.com/pybind/pybind11>

⁶ <https://github.com/kimiyoung/planetoid>

表4: GRAND+的超参数配置。

		$\tilde{\alpha}$	wr	L	$d_{mh,max}$	k	N	λ	m	a	x
Cora	GRAND+ (P)	10^{-2}	10^{-3}	2	64	10^{-7}	32	20	1.5		
	GRAND+ (A)	10^{-2}	10^{-3}	2	64	10^{-7}	32	4	1.5		
	GRAND+ (S)	10^{-2}	10^{-3}	2	64	10^{-7}	32	2	1.5		
馨香坊	盛大+ (P)	10^{-3}	10^{-3}	2	256	10^{-7}	32	10	0.8		
	盛大+ (a)	10^{-3}	10^{-3}	2	256	10^{-7}	32	2	0.8		
	grand+ (s)	10^{-3}	10^{-3}	2	256	10^{-7}	32	2	0.8		
公共医学	盛大+ (P)	10^{-2}	10^{-2}	1	-	10^{-5}	16	6	1.0		
	盛大+ (a)	10^{-2}	10^{-2}	1	-	10^{-5}	16	4	1.0		
	盛大+ (s)	10^{-2}	10^{-2}	1	-	10^{-5}	16	2	1.0		
AMiner-CS	盛大+ (P)	10^{-2}	10^{-2}	1	-	10^{-5}	64	6	1.5		
	盛大+ (a)	10^{-2}	10^{-2}	1	-	10^{-5}	64	4	1.5		
	盛大+ (s)	10^{-2}	10^{-2}	1	-	10^{-5}	64	2	1.5		
瑞德	grand+ (p)	10^{-4}	0	2	512	10^{-5}	64	6	1.5		
	grand+ (a)	10^{-4}	0	2	512	10^{-5}	64	6	1.5		
	盛大+ (s)	10^{-4}	0	2	512	10^{-7}	64	2	1.5		
Amazon2M	大+ (p)	10^{-3}	10^{-5}	2	1024	10^{-6}	64	6	0.8		
	盛大+ (a)	10^{-3}	10^{-5}	2	1024	10^{-6}	64	4	0.8		
	大+ (s)	10^{-3}	10^{-5}	2	1024	10^{-6}	32	2	0.8		
MAG-Scholar-C	Grand+ (p)	10^{-2}	0	2	32	10^{-5}	32	10	1.0		
	GRAND+ (A)	10^{-2}	0	2	32	10^{-5}	32	10	1.0		
	盛大+ (s)	10^{-2}	0	2	32	10^{-5}	32	2	1.0		

GRAND+的超参数选择包括两个阶段: 我们首先对神经网络的基本超参数进行搜索。具体来说, 我们搜索学习率 $\square\square$, 从 $\{10^{-2}, 10^{-3}, 10^{-4}\}$, 重量衰减率 wr 来自 $\{0, 10^{-5}, 10^{-3}, 10^{-2}\}$, 隐层的数量 L_m 来自 $\{1, 2\}$, 隐层的尺寸 d_h 来自 $\{32, 64, 128, 256, 512, 1024\}$ 。

在第二阶段, 我们将这些基本超参数固定为最佳配置, 并搜索以下具体的超参数: DropNode rate \square , augmentation times per batch \square , threshold r_{max} , max neighborhood size \square , propagation order N 、信度阈值 γ , 最大一致性损失权重 λ_{max} , 未标记子集的大小 $|\square|$ 和一致性损失函数 D 。为了降低搜索成本, 我们将一些超参数固定下来。具体来说, 我们在所有的数据集上固定 $\delta=0.5$, $M=2$ 和 $\gamma=2$ 。对于Cora、Pubmed和Citeseer, 我们设定 $|\square|=|\square|$, 而对于其他数据集, 我们设定 $|\square|=10000$ 。我们还提供了一个分析, 以了解以下因素的影响[11], 见附录A.3。我们采用KL发散作为一致性。

对AMiner-CS、Reddit和Amazon2M使用损失函数, 而对其他数据集使用 \square_2 距离。这是因为 \square_2 距离在处理有大量类的数据集时容易出现梯度消失的问题。然后, 我们对 $\square\square\square\square$ 、 \square 、 \square 和 $\square\square\square\square$ 进行超参数选择。具体而言, 我们从以下方面搜索 $\square\square\square\square$ $\{10^{-5}, 10^{-6}, 10^{-7}\}$, k 来自 $\{16, 32, 64, 128\}$, N 来自 $\{2, 4, 6, 8, 10, 20\}$ 和 $\square\square\square\square\{0.5, 0.8, 1.0, 1.2, 1.5\}$ 。表4列出了GRAND+的最佳超参数配置。

A.2 定理证明

为了证明定理1, 我们首先介绍以下定理:

