

基于参数差异假设的图卷积网络对抗性攻击

吴翼腾¹, 刘 伟¹, 于淑乔²

(1. 信息工程大学信息技术研究所, 河南郑州 450002; 2. 墨尔本大学, 澳大利亚墨尔本 3010)

摘要: 神经网络容易受到对抗性攻击安全威胁. 现有神经网络对抗性攻击思想可以概括为构造矛盾的训练数据. 矛盾数据假设不能很好地解释神经网络过拟合训练数据的攻击场景. 本文以有效攻击前后神经网络模型的训练参数应该具有较大差异为基本出发点, 以图卷积网络为具体研究对象, 建立基于参数差异假设的对抗性攻击模型. 将统计诊断的重要结果Cook距离引入对抗性攻击, 提出基于Cook距离的参数差异度量方法. 采用基于Cook距离梯度的攻击方法, 首次得出了攻击梯度的闭式解, 并结合梯度下降算法思想和贪心算法思想提出完整的攻击算法. 最后设计实验验证了参数差异假设的合理性和基于该假设导出方法的有效性; 验证了梯度信息对图场景离散数据的可用性; 仿真示例说明了攻击梯度闭式解的正确性; 与其他攻击方法对比分析了攻击方法的有效性.

关键词: 图卷积网络; 对抗性攻击; 矛盾数据假设; 参数差异假设; Cook距离

基金项目: 自然科学基金创新研究群体项目(No.61521003); 国家重点研发计划(No.2016QY03D0502); 郑州市协同创新重大专项基金(No.162/32410218)

中图分类号: TP18

文献标识码: A

文章编号: 0372-2112(2023)02-0330-12

电子学报 URL: <http://www.ejournal.org.cn>

DOI: 10.12263/DZXB.20210222

Adversarial Attacks on Graph Convolution Networks Based on Parameter Discrepancy Hypothesis

WU Yi-teng¹, LIU Wei¹, YU Xu-qiao²

(1. Institute of Information Technology, Information Engineering University, Zhengzhou, Henan 450002, China;

2. University of Melbourne, Melbourne 3010, Australia)

Abstract: Graph neural networks (GNNs) are vulnerable to adversarial attacks. Existing GNN adversarial attacks can be generalized as constructing contradictory training data. However, the existing methods based on contradictory data hypothesis cannot explain well why the false outputs could be generated when GNNs fit the training data well. Firstly, based on the discrepancy of model parameters of GNNs before and after attack, poisoning attack model is proposed taking graph convolution network as a target. Secondly, a parameter difference metric, Cook distance, is proposed. The closed form solution of attack gradients is obtained for the first time, and an attack algorithm is given based on the idea of gradient descent and greedy algorithm. Finally, the rationality of the hypothesis of parameter discrepancy and the effectiveness of the proposed method are verified by experiments; the availability of gradients to discrete data of graph is verified; the correctness of closed form solution of attack gradients is illustrated by a numerical example; the effectiveness of attack method is analyzed compared with other attacks.

Key words: graph convolutional network; adversarial attack; contradictory data hypothesis; parameter discrepancy hypothesis; Cook distance

Foundation Item(s): National Natural Science Foundation of China (No.61521003); National Key R&D Plan (No.2016QY03D0502); Major Special Fund for Zhengzhou Collaborative Innovation (No.162/32410218)

1 引言

近年来随着深度学习方法的广泛应用, 其安全问题也受到普遍关注. 2014年Szegedy等人^[1]在研究图像数据的卷积神经网络(Convolutional Neural Networks,

CNN)模型安全问题时提出“对抗样本”的概念. 图像中的对抗样本至今已从实验室水平发展至攻防实战阶段^[2]. 研究对抗性攻击方法和对抗样本的产生机理^[3]是深度学习方法应用于无人驾驶、智能安防等安全关键

性领域的必要前提^[4].

图数据与语义信息相结合具有表达现实数据的天然优势. 传统的复杂网络研究主要聚焦于节点和边的连接模式, 研究图中关于节点、连边和局部结构的性质^[5], 难以将图数据和对应的属性信息有机结合. 深度学习方法在图像、文本等领域取得了目前已知的最好结果, 却难以直接处理图数据. 为有效处理图数据和语义属性数据, 诞生了以图卷积网络^[6-9]为代表的图神经网络^[10]. 图神经网络是专门针对图数据设计的端到端的深度学习模型^[11-13], 为含有图数据的应用场景提供了极具竞争力的学习方案.

与其他深度学习技术类似, 图神经网络同样存在安全威胁. 有目的地对网络结构或属性特征施加微小扰动, 会使图神经网络在具体任务中高置信度地给出错误输出. 2018年, Zügner等人^[14]首次提出图神经网络对抗性攻击. 该课题近两年来逐渐受到更多关注^[15-20]. 按污染数据的阶段划分, 对抗性攻击分为污染训练数据的投毒攻击和污染测试数据的逃逸攻击. 与以往对图神经网络投毒攻击的理解不同, 本文将图投毒攻击归结为攻击前后模型参数的显著差异. 现有文献^[14, 18, 19, 21, 22]将投毒攻击形式化表达为, 寻求攻击方法, 使重训练的图神经网络在训练集上的损失函数达到最大, 以降低测试数据上相关任务的准确率. 该表述所描述的投毒攻击可以概括为: 寻求攻击方法, 构造一组存在矛盾的训练数据, 使图神经网络难以对其拟合, 最终学习出一个“坏模型”, 使其在测试集的效果显著下降. 矛盾数据假设是对投毒攻击本质的一种重要概括, 但并不全面. 例如, 存在这样的攻击场景——构造投毒的训练数据, 使图神经网络能够对其很好地拟合, 却难以对测试数据正确分类, 即最终学习出一个“假模型”. 分析针对测试数据的实际攻击场景, 模型结构和测试数据不变, 发生改变的只有模型训练参数, 可见模型训练参数对预测结果有较大影响. 本文将以此为动机展开研究.

基于以上分析, 本文以攻击前后图神经网络的参数差异为基本出发点, 建立基于参数差异假设的对抗性攻击模型. 主要面临如下挑战:

(1) 如何衡量参数差异. 这需要建立参数差异的有效度量.

(2) 如何求解攻击模型. 由于图神经网络处理的数据为离散数据, 而传统观点^[21]认为基于梯度的方法不适用于离散数据, 但是在图神经网络攻击这一特殊场景中, 基于梯度的方法是否有效不能一概而论, 需要具体分析.

(3) 若基于梯度的方法有效, 能否得出攻击梯度的闭式解. 传统使用反向传播算法求解攻击梯度, 无法得出攻击梯度的显式表达式, 使得攻击方法具有黑盒性, 不利于解析攻击原理和进一步寻找攻击问题的本质.

本文的主要贡献如下:

(1) 基于“有效攻击前后图神经网络训练参数应该具有较大差异”的假设^[17](简称参数差异假设)实施攻击, 以图卷积网络为研究对象, 建立新的对抗性攻击模型.

(2) 从统计诊断经典文献^[23~29]中引入Cook距离用以衡量攻击前后的参数差异. 设计实验验证参数差异度量方法的有效性、参数差异假设的合理性和其导出方法的有效性.

(3) 推导出Cook距离对邻接矩阵和特征矩阵攻击梯度的闭式解, 并对理论结果进行仿真. 设计实验验证图卷积网络对抗性攻击情景中离散数据基于梯度算法的有效性和可用性. 基于攻击梯度, 根据梯度下降法和贪心算法的思想设计攻击算法, 并进行对比实验.

2 基本概念

2.1 图和图卷积网络

图表示为 $G(V, E)$, 其中 V 表示节点集合, E 表示连边集合. 设节点数 $|V|=N$, 则无权无向图可用对称的邻接矩阵 $A=\{0, 1\}^{N \times N}$ 表示, $A^T=A$. 图中每个节点有 n 维的特征向量, 节点特征用矩阵 $X=\{0, 1\}^{N \times n}$ 或 $X \in \mathbb{R}^{N \times n}$ 表示. 本文使用的特征矩阵为离散情况, $X=\{0, 1\}^{N \times n}$. 文献^[7~9]将图卷积网络简化为SGC(Simple Graph Convolution), 它具有低通滤波器的作用. 本文以图卷积网络SGC为攻击和研究对象. SGC模型的公式表达式为

$$Y_{\text{out}} = \text{softmax}(AXW) \quad (1)$$

其中, Y_{out} 为模型输出; A 为滤波矩阵; X 为输入特征向量; W 为参数矩阵. 在文献^[7]中, A 的形式通常为

$$A = \tilde{L}^k, \tilde{L} = D^{-\frac{1}{2}} \tilde{A} D^{-\frac{1}{2}}, \tilde{A} = I + A, \tilde{D} = \text{diag}(\tilde{A} \mathbf{1}) \quad (2)$$

记 $\text{softmax}(\cdot)$ 为 $\sigma(\cdot)$, Y 表示 one-hot 编码的标签矩阵. $\mathbf{1}$ 表示 N 维全 1 列向量. 使用交叉熵损失函数, 可得到矩阵形式的损失函数计算公式为

$$\mathcal{L} = -\text{tr}[Y^T \ln[\sigma(ZW)]] \quad (3)$$

2.2 图卷积网络的对抗性攻击

本文研究非指定目标、数据投毒攻击. 非指定目标攻击不指定具体的一个或几个攻击目标, 需要使测试集的准确率整体下降; 投毒攻击指允许图卷积网络对投毒的训练数据重新训练, 重训练的图卷积网络在测试集的准确率仍然下降. 文献^[14, 18, 19, 21, 22]建立了图卷积网络的对抗样本生成或投毒攻击模型. 按上述文献定义, 攻击方法可以统一概括为以下约束优化问题:

$$(\hat{A}, \hat{X}) = \arg \max_{(\hat{A}, \hat{X})} \mathcal{L}(W^*; \hat{A}, \hat{X}) \quad (4)$$

$$\begin{aligned} \text{s.t. } W^* &= \arg \min_W \mathcal{L}(W; \hat{A}, \hat{X}) \\ \|\hat{A} - A\|_0 &+ \|\hat{X} - X\|_0 \leq \delta \end{aligned} \quad (5)$$

其中, \mathbf{A} 和 \mathbf{X} 是原邻接矩阵和特征矩阵; $\hat{\mathbf{A}}, \hat{\mathbf{X}}$ 是扰动后的邻接矩阵和特征矩阵; $\|\cdot\|_0$ 为矩阵中非 0 元素的个数; δ 是扰动开销; \mathbf{W}^* 是扰动后得到的训练参数. 投毒攻击允许对参数重新训练, 属于双层优化问题. 若损失函数表示指定节点的损失函数, 为指定目标攻击; 若为全局损失函数, 为非指定目标攻击. 文献[22]研究了非指定目标数据投毒攻击.

从上述投毒攻击的形式化表述中可以看出, 现有攻击方法都是构造扰动后的样本数据, 使得图卷积网络在扰动后数据集上损失函数最大, 从而训练出不良参数. 现有方法可概括为基于“有效攻击的训练集中存在矛盾的训练数据”假设(简称矛盾数据假设)的攻击方法, 攻击的实质是构造矛盾的训练数据.

3 基于参数差异假设的攻击方法

3.1 基于参数差异假设的投毒攻击模型

基于“有效攻击前后图卷积网络模型参数应该具有较大差异”的参数差异假设, 建立图卷积网络的投毒攻击模型. 基于参数差异假设的投毒攻击模型可表示为如下约束优化问题:

$$\begin{aligned} (\hat{\mathbf{A}}, \hat{\mathbf{X}}) &= \arg \max_{(\hat{\mathbf{A}}, \hat{\mathbf{X}})} \|\mathbf{W}^* - \mathbf{W}_0\|_M^2 & (6) \\ \text{s.t. } \mathbf{W}^* &= \arg \min_{\mathbf{W}} \mathcal{L}(\mathbf{W}; \hat{\mathbf{A}}, \hat{\mathbf{X}}) & (7) \\ \|\hat{\mathbf{A}} - \mathbf{A}\|_0 + \|\hat{\mathbf{X}} - \mathbf{X}\|_0 &\leq \delta \end{aligned}$$

其中, $\mathbf{W}_0 = \arg \min_{\mathbf{W}} \mathcal{L}(\mathbf{W}; \mathbf{A}, \mathbf{X})$ 为图卷积网络的最佳训练参数; $\|\cdot\|_M$ 是矩阵的某种加权范数, 用以衡量参数差异, 即

$$\|\mathbf{W}^* - \mathbf{W}_0\|_M = \sqrt{\text{vec}^T(\mathbf{W}^* - \mathbf{W}_0) \mathbf{M} \text{vec}(\mathbf{W}^* - \mathbf{W}_0)} \quad (8)$$

其中, \mathbf{M} 是 n 阶实对称权矩阵; \mathbf{W}^* 是邻接矩阵或特征矩阵扰动后图卷积网络的训练参数; $\text{vec}(\cdot)$ 是矩阵按列优先的拉直运算. 设模型参数采用梯度下降法训练, 即

$$\mathbf{W}^t = \mathbf{W}^{t-1} - \alpha \cdot \nabla_{\mathbf{W}} \mathcal{L}(\mathbf{W}^{t-1}) \quad (9)$$

则经过 $t=t_0$ 轮训练, 可得模型的训练参数 $\mathbf{W}^* = \mathbf{W}^{t_0}$. 具体地, 对于式(1)的图卷积网络和式(3)的损失函数, 可求得

$$\nabla_{\mathbf{W}} \mathcal{L}(\mathbf{W}) = (\hat{\mathbf{A}} \hat{\mathbf{X}})^T \left[\sigma(\hat{\mathbf{A}} \hat{\mathbf{X}} \mathbf{W}) - \mathbf{Y} \right] \quad (10)$$

$\hat{\mathbf{A}}$ 和 $\hat{\mathbf{X}}$ 表示扰动后的 \mathbf{A} 和扰动后的 \mathbf{X} , 设 $\hat{\mathbf{Z}} = \hat{\mathbf{A}} \hat{\mathbf{X}}$, 则有

$$\mathbf{W}^* = \mathbf{W}^0 - \alpha \cdot \sum_{t=0}^{t_0-1} \hat{\mathbf{Z}}^T \left[\sigma(\hat{\mathbf{Z}} \mathbf{W}^t) - \mathbf{Y} \right] \quad (11)$$

注意, 投毒攻击场景下, 图卷积网络正向训练时所

用的数据为扰动后的邻接矩阵 $\hat{\mathbf{A}}$ 和特征矩阵 $\hat{\mathbf{X}}$.

3.2 衡量参数差异的 Cook 距离及其计算方法

设攻击后模型求得的参数 \mathbf{W} 的估计值为 \mathbf{W}^* , 原始训练参数为 \mathbf{W}_0 . 首先建立某种“广义距离”对二者差异进行度量. 本文引入统计诊断中的重要成果 Cook 距离及其近似方法^[23-25, 27-29] 衡量参数差异.

定义 1 (Cook 距离) 参数矩阵 \mathbf{W}^* 与 \mathbf{W}_0 的 Cook 距离 CD 定义为

$$\text{CD} = \text{vec}^T \left[\nabla_{\mathbf{W}} \mathcal{L}(\mathbf{W}^*) \right] \left[\nabla_{\mathbf{W}}^2 \mathcal{L}(\mathbf{W}^*) \right]^{K-2} \text{vec} \left[\nabla_{\mathbf{W}} \mathcal{L}(\mathbf{W}^*) \right] \quad (12)$$

其中, $\left[\nabla_{\mathbf{W}}^2 \mathcal{L}(\mathbf{W}^*) \right]^{K-2}$ 是权矩阵, K 为整数, 通常 $K=1, 2, 3$.

注意, 若 CD 严格表示距离, 本该满足距离的公理化定义, 但这里因循统计诊断之约定, Cook 距离仍按式(12)定义. 式(12)中 $\nabla_{\mathbf{W}} \mathcal{L}(\mathbf{W}^*) = \mathbf{Z}^T [\sigma(\mathbf{Z} \mathbf{W}^*) - \mathbf{Y}]$, $\mathbf{Z} = \mathbf{A} \mathbf{X}$, 即区别于式(10)和式(11), 此处代入未扰动的邻接矩阵和特征矩阵. 这也是本文基于参数差异假设建立攻击模型与基于矛盾数据假设建立攻击模型的重要区别. 第 4.3 节和第 4.4 节将设计实验验证参数差异假设的合理性和有效性, 并与矛盾数据假设做对比分析.

Cook 距离 CD 的具体计算式由定理 1 给出.

定理 1 (Cook 距离的计算公式) 式(12)的 Cook 距离 CD 可写成如下形式:

$$\begin{aligned} \text{CD} &= \text{vec}^T \left(\mathbf{Z}^T \mathbf{E} \right) \left[\left(\mathbf{I}_m \otimes \mathbf{Z}^T \right) \mathbf{K}_{mN} \mathbf{A}^* \mathbf{K}_{Nm} \left(\mathbf{I}_m \otimes \mathbf{Z} \right) \right]^{K-2} \\ &\quad \times \text{vec} \left(\mathbf{Z}^T \mathbf{E} \right) \end{aligned} \quad (13)$$

其中, Hessian 矩阵

$$\begin{aligned} \mathbf{M}(\mathbf{W}^*) &= \nabla_{\mathbf{W}}^2 \mathcal{L}(\mathbf{W}^*) \\ &= \left(\mathbf{I}_m \otimes \mathbf{Z}^T \right) \mathbf{K}_{mN} \mathbf{A}^* \mathbf{K}_{Nm} \left(\mathbf{I}_m \otimes \mathbf{Z} \right) \\ \mathbf{E} &= \mathbf{Y} - \sigma(\mathbf{Z} \mathbf{W}^*) \end{aligned}$$

为残差矩阵. “ \otimes ”表示矩阵的 Kronecker 积, 有

$$\mathbf{A} \otimes_{m \times n} \mathbf{B} \otimes_{p \times q} = [\mathbf{A}_{ij} \mathbf{B}]_{mp \times nq} \quad (14)$$

\mathbf{A}^* 为块对角阵, 有

$$\begin{aligned} \mathbf{A}^* &= \begin{bmatrix} \lambda_1 & & & \\ & \ddots & & \\ & & \ddots & \\ & & & \lambda_N \end{bmatrix} \\ \lambda_i &= \frac{\partial \sigma(\xi_i)}{\partial \xi_i} \end{aligned} \quad (15)$$

$$\mathbf{Z} \mathbf{W}^* = \begin{bmatrix} \xi_1^T \\ \vdots \\ \xi_N^T \end{bmatrix}$$

$$\begin{aligned}\lambda_i &= \frac{\partial \sigma(\xi_i)}{\partial \xi_i} \\ &= \frac{\partial \text{softmax}(\xi_i)}{\partial \xi_i} \\ &= \frac{\mathbf{1}^T \exp(\xi_i) \text{diag}[\exp(\xi_i)] - \exp(\xi_i) \exp^T(\xi_i)}{[\mathbf{1}^T \exp(\xi_i)]^2}\end{aligned}\quad (16)$$

定理 1 的证明将在附录中给出. 当 $K=1$ 时, 涉及 Hessian 矩阵 $\nabla_{\mathbf{W}}^2 \mathcal{L}(\mathbf{W}^*)$ 求逆, 实用中 $\nabla_{\mathbf{W}}^2 \mathcal{L}(\mathbf{W}^*)$ 通常不可逆, 因此 $K=1$ 时仅用于理论分析. 实际计算 CD 时为避免求逆而取 $K=2, 3$, 或在具体场景中寻找 $K=1$ 的等价方法或近似方法. 本文取 $K=2, 3$.

3.3 攻击梯度的闭式解

对于邻接矩阵 $\mathbf{A}=\{0, 1\}^{N \times N}$ 和特征矩阵 $\mathbf{X}=\{0, 1\}^{N \times n}$ 均为离散值的情况, 显然优化问题(式(6)、式(7))属于 NP 难问题. 可以使用动态规划、贪心算法等求解. 以贪心算法为例, 在扰动开销 δ 的范围内, 每增加第 i 个扰动元素, 对特征矩阵 \mathbf{X} 的攻击需要计算 $Nn-i+1$ 次扰动对 CD 的影响, 即大约需要计算 δNn 次 CD (对邻接矩阵 \mathbf{A} 扰动大约需要计算 δN^2 次 CD). 为减少计算 CD 的次数, 基于参数差异的攻击模型设计了利用 CD 梯度信息的求解算法. 对于具体的图卷积网络, 式(13)建立了 CD 关于扰动后邻接矩阵 $\hat{\mathbf{A}}$ 和扰动后特征矩阵 $\hat{\mathbf{X}}$ (包含在参数 \mathbf{W}^* 中) 的显式表达式, 根据该表达式可以推导出 CD 对矩阵 $\hat{\mathbf{A}}$ 和 $\hat{\mathbf{X}}$ 梯度的闭式解. 闭式解利于理论研究、解析图卷积网络的攻击原理, 容易对产生攻击的主要影响因素追踪定位和溯源, 以避免反向传播等梯度求解方法的黑盒性.

采用基于梯度的攻击方法, 首先求解 Cook 距离 CD 对扰动后的邻接矩阵和特征矩阵的梯度, 可得关于攻击梯度的如下两个定理.

定理 2 (特征矩阵的攻击梯度) 设图卷积网络(式(1))和损失函数(式(3))通过梯度下降法(式(9))训练, 将权矩阵 $\mathbf{M}^K(\mathbf{W}^*)$ 与式(9)中的参数矩阵 \mathbf{W}^t ($t=0, 1, 2, \dots, t_0-1$) 视为常数矩阵. 则 Cook 距离 CD 对扰动后的特征矩阵 $\hat{\mathbf{X}}$ 的梯度可以表示为

$$\nabla_{\hat{\mathbf{X}}} \text{CD} = 2\alpha \cdot \hat{\mathbf{A}} \sum_{t=0}^{t_0-1} (\mathbf{Q}^t \mathbf{W}^{tT} - \hat{\mathbf{E}}^t \mathbf{H}_K^T) \quad (17)$$

其中, $\hat{\mathbf{E}}^t = \mathbf{Y} - \sigma(\hat{\mathbf{Z}} \mathbf{W}^t)$ 为梯度下降法训练至第 t 轮 (即 $\mathbf{W} = \mathbf{W}^t$) 时的残差矩阵, $\mathbf{H}_1 = \mathbf{Z}^T \mathbf{E}$, 即

$$\text{vec}(\mathbf{H}_K) = \mathbf{M}(\mathbf{W}^*) \text{vec}(\mathbf{H}_{K-1}) \quad (18)$$

$$\text{vec}(\mathbf{Q}^t) = \hat{\mathbf{A}}^t \text{vec}(\mathbf{H}_K^T \hat{\mathbf{Z}}^T) \quad (19)$$

$\hat{\mathbf{A}}^t$ 的形式同式(15), 只是将其中的 $\mathbf{Z} \mathbf{W}^*$ 换成 $\hat{\mathbf{Z}} \mathbf{W}^*$.

定理 2 的证明将在附录中给出.

定理 3 (邻接矩阵的攻击梯度) 设图卷积网络(式(1))和损失函数(式(3))通过梯度下降法(式(9))训练, 将权矩阵 $\mathbf{M}^K(\mathbf{W}^*)$ 、含自连边的度矩阵 $\tilde{\mathbf{D}}$ 与式(11)中的参数矩阵 \mathbf{W}^t ($t=0, 1, 2, \dots, t_0-1$) 视为常数矩阵. 则 Cook 距离 CD 对扰动后的含自连边的邻接矩阵 $\hat{\mathbf{A}}$ 的梯度可以表示为

$$\nabla_{\hat{\mathbf{A}}} \text{CD} = 2\alpha \tilde{\mathbf{D}}^{-\frac{1}{2}} \left[\sum_{i=0}^{k-1} \hat{\mathbf{L}}^i \left(\sum_{t=0}^{t_0-1} \mathbf{Q}^t \mathbf{W}^{tT} - \hat{\mathbf{E}}^t \mathbf{H}_K^T \right) \left(\hat{\mathbf{L}}^{k-i-1} \hat{\mathbf{X}} \right)^T \right] \tilde{\mathbf{D}}^{-\frac{1}{2}} \quad (20)$$

其中, $\hat{\mathbf{L}} = \tilde{\mathbf{D}}^{-\frac{1}{2}} \hat{\mathbf{A}} \tilde{\mathbf{D}}^{-\frac{1}{2}}$. 特别地, 当 $k=1$ 时, 有

$$\nabla_{\hat{\mathbf{A}}} \text{CD} = 2\alpha \tilde{\mathbf{D}}^{-\frac{1}{2}} \left(\sum_{t=0}^{t_0-1} \mathbf{Q}^t \mathbf{W}^{tT} - \hat{\mathbf{E}}^t \mathbf{H}_K^T \right) \hat{\mathbf{X}}^T \tilde{\mathbf{D}}^{-\frac{1}{2}} \quad (21)$$

定理 3 的证明将在附录中给出. 注意, CD 的计算式(13)中, $\mathbf{Z} = \mathcal{A} \mathbf{X}$, 其中 \mathcal{A} 和 \mathbf{X} 为未扰动的归一化邻接矩阵和特征矩阵而非扰动后的矩阵, 因此求导时看作常数矩阵. 扰动信息体现在参数 \mathbf{W}^* 的训练过程(式(11))之中. 这也体现了 Cook 距离 CD 的物理意义, 即扰动后的数据诱导图卷积网络训练出一组异于原参数的参数 \mathbf{W}^* , 使新参数 \mathbf{W}^* 偏离未扰动的损失函数 \mathcal{L} 的驻点, 达到攻击目的.

3.4 基于参数差异假设的投毒攻击算法

根据上述理论, 基于梯度下降算法与贪心算法的思想设计基于参数差异假设的图卷积网络投毒攻击算法 (Parattack). 攻击邻接矩阵或特征矩阵算法的主要步骤包括:

(1) 根据上一轮的扰动列表 `disturb_list` 求得扰动后的邻接矩阵 $\hat{\mathbf{A}}$ 和特征矩阵 $\hat{\mathbf{X}}$ 为 $\hat{\mathbf{A}}, \hat{\mathbf{X}} = \text{get_modified}(\mathbf{A}, \mathbf{X}, \text{disturb_list})$.

(2) 图卷积网络进行正向训练, 得训练参数 $\mathbf{W}^* = \text{train_GNN}(\hat{\mathbf{A}}, \hat{\mathbf{X}}, \mathbf{Y})$.

(3) 根据式(17)和式(20)或其等价方法求攻击梯度矩阵 $\nabla_{\hat{\mathbf{A}}} \text{CD}, \nabla_{\hat{\mathbf{X}}} \text{CD}$ 为 $\nabla_{\hat{\mathbf{A}}} \text{CD}, \nabla_{\hat{\mathbf{X}}} \text{CD} = \text{get_CD_gradient}(\hat{\mathbf{A}}, \hat{\mathbf{X}}, \mathbf{W}^*, \mathbf{Y})$.

(4) 对攻击梯度矩阵 $\nabla_{\hat{\mathbf{A}}} \text{CD}, \nabla_{\hat{\mathbf{X}}} \text{CD}$ 的符号进行处理, 保留对扰动筛选有效的攻击梯度, 得许用梯度矩阵 $\mathbf{G}_{\hat{\mathbf{A}}}, \mathbf{G}_{\hat{\mathbf{X}}}, \mathbf{G}_{\hat{\mathbf{A}}}, \mathbf{G}_{\hat{\mathbf{X}}} = \text{sign_process}(\nabla_{\hat{\mathbf{A}}} \text{CD}, \nabla_{\hat{\mathbf{X}}} \text{CD}, \hat{\mathbf{A}}, \hat{\mathbf{X}})$. 对于扰动后的邻接矩阵 $\hat{\mathbf{A}}$ 或特征矩阵 $\hat{\mathbf{X}}$ 的元素值为 1 的位置, 当梯度值为负数时有效. 对于元素值为 0 的位置, 当梯度值为正数时有效; 为保证无权无向网络邻接矩阵的对称性, 再将邻接矩阵的有效梯度转置相加, 得 $\mathbf{G}_{\hat{\mathbf{A}}}, \mathbf{G}_{\hat{\mathbf{X}}}, \mathbf{G}_{\hat{\mathbf{A}}}, \mathbf{G}_{\hat{\mathbf{X}}}$ 称为许用梯度矩阵. 许用梯度矩阵的具体计算方式如下: 首先计算 $\mathbf{S}_{\hat{\mathbf{A}}} = (-2\hat{\mathbf{A}} + \mathbf{1}) \odot \nabla_{\hat{\mathbf{A}}} \text{CD}$, $\mathbf{S}_{\hat{\mathbf{X}}} = (-2\hat{\mathbf{X}} + \mathbf{1}) \odot \nabla_{\hat{\mathbf{X}}} \text{CD}$; 然后将 $\mathbf{S}_{\hat{\mathbf{A}}}, \mathbf{S}_{\hat{\mathbf{X}}}$ 中的正值置为 1, 负值置

为0,得许用位置矩阵 F_A, F_X ;最后求得许用梯度矩阵 $G_A = F_A \odot S_A + (F_A \odot S_A)^T, G_X = F_X \odot S_X$.

(5)根据许用梯度矩阵和攻击类型(攻击邻接矩阵、特征矩阵或邻接矩阵和特征矩阵的联合攻击),每轮筛选许用梯度的绝对值最大的 $\gamma n_0 (0 < \gamma < 1)$ 个位置,对 \hat{A}, \hat{X} 增量式地进行扰动,达到扰动总数 n_0 ;达到扰动总数后,继续对扰动列表进行更新,直至满足迭代总次数iters停止. $\text{disturb_list} = \text{disturb_AX}(G_A, G_X, \hat{A}, \hat{X}, n_0)$.

算法伪代码如表算法1所示.

算法1 基于参数差异假设的图卷积网络投毒攻击算法(Parattack)

输入: 邻接矩阵 A ,特征矩阵 X ,标签 Y ,攻击点数 n_0 ,迭代次数iters.

输出: 扰动列表disturb_list.

disturb_list=[];

FOR i in range (iters):

$\hat{A}, \hat{X} = \text{get_modified}(A, X, \text{disturb_list});$

$W^* = \text{train_GNN}(\hat{A}, \hat{X}, Y);$

$\nabla_{\hat{A}} \text{CD}, \nabla_{\hat{X}} \text{CD} = \text{get_CD_gradient}(\hat{A}, \hat{X}, W^*, Y);$

$G_A, G_X = \text{sign_process}(\nabla_{\hat{A}} \text{CD}, \nabla_{\hat{X}} \text{CD}, \hat{A}, \hat{X});$

$\text{disturb_list} = \text{disturb_AX}(G_A, G_X, \hat{A}, \hat{X}, n_0);$

RETURN disturb_list

该算法是基于梯度下降算法和贪心算法的思想设计的,当 $\gamma = 1/n_0$,每次允许不重复地扰动1个元素时,即为贪心算法.然而,图卷积网络对抗攻击中的优化对象为邻接矩阵和特征矩阵,二者中的元素均为离散的0,1整数变量,与连续取值的情况不同.从理论上,我们暂未得到算法收敛性的有关结论.尽管暂未从理论上证明该算法的收敛性,算法仍然具有有效性和可用性.本文通过实验分析的方法,设计实验说明算法的可用性,即攻击梯度对本文场景的离散数据扰动筛选仍然具有重要的指导意义.

4 实验验证

实验验证部分首先对攻击梯度的闭式解给出算例;其次设计实验说明基于梯度的算法对于图卷积网络对抗攻击场景的可用性;然后针对节点分类任务,设计实验验证参数度量的合理性、参数差异假设的合理性和其导出方法的有效性;最后将本文提出的基于参数差异假设的攻击方法与其他方法进行对比实验.实验中采用4个数据集,即polblogs^[30]数据集、cora数据集、cora_ml^[31]数据集和citeseer^[32]数据集.各个数据集的统计特性如表1所示.

4.1 攻击梯度闭式解的算例

设邻接矩阵 A 、特征矩阵 X 、初始参数 W^0 和标签矩阵 Y 分别为

表1 实验数据集的统计特性

数据集	节点数	连边数	特征维数	分类数
polblogs	1 222	16 714	1 490	2
cora	2 485	5 069	1 433	7
cora_ml	2 810	7 981	2 879	7
citeseer	2 110	3 668	3 703	6

$$A = \begin{bmatrix} 0 & 1 & 1 & 0 \\ 1 & 0 & 1 & 1 \\ 1 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{bmatrix}, X = \begin{bmatrix} 1 & 0 & 1 \\ 0 & 1 & 1 \\ 1 & 0 & 1 \\ 1 & 1 & 1 \end{bmatrix}$$

$$W^0 = \begin{bmatrix} 1.2 & 1.9 & 2.1 \\ 1.8 & 1.1 & 1.3 \\ 2.0 & 2.2 & 2.1 \end{bmatrix}, Y = \begin{bmatrix} 0 & 0 & 1 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \end{bmatrix}$$

设扰动后的邻接矩阵 \hat{A} 和特征矩阵 \hat{X} 分别为

$$\hat{A} = \begin{bmatrix} 0 & 1 & 1 & 1 \\ 1 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 \end{bmatrix}, \hat{X} = \begin{bmatrix} 1 & 0 & 1 \\ 0 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 0 \end{bmatrix}$$

设图卷积网络(式(1))使用梯度下降法(式(9))训练的迭代次数为 $t_0 = 3$,学习率 $\alpha = 0.1$.经过3轮训练得到的参数和矩阵 H_3 为

$$W^1 = \begin{bmatrix} 1.30 & 1.79 & 2.12 \\ 1.90 & 1.00 & 1.31 \\ 2.04 & 2.15 & 2.11 \end{bmatrix}$$

$$W^2 = \begin{bmatrix} 1.37 & 1.70 & 2.13 \\ 1.97 & 0.92 & 1.31 \\ 2.08 & 2.10 & 2.12 \end{bmatrix}$$

$$W^3 = \begin{bmatrix} 1.43 & 1.62 & 2.15 \\ 2.04 & 0.85 & 1.31 \\ 2.10 & 2.07 & 2.13 \end{bmatrix}$$

$$H_3 = \begin{bmatrix} 1.75 & -1.86 & 0.10 \\ 1.26 & -1.31 & 0.05 \\ 2.48 & -2.62 & 0.14 \end{bmatrix}$$

$\hat{E}^t (t=0, 1, 2)$ 和 E 矩阵为

$$\hat{E}^0 = \begin{bmatrix} -0.25 & -0.31 & 0.56 \\ 0.72 & -0.31 & -0.41 \\ -0.28 & -0.33 & 0.61 \\ 0.74 & -0.30 & -0.44 \end{bmatrix}$$

$$\hat{E}^1 = \begin{bmatrix} -0.31 & -0.24 & 0.55 \\ 0.67 & -0.26 & -0.41 \\ -0.32 & -0.29 & 0.60 \\ 0.68 & -0.23 & -0.45 \end{bmatrix}$$

$$\hat{E}^2 = \begin{bmatrix} -0.36 & -0.19 & 0.55 \\ 0.63 & -0.22 & -0.41 \\ -0.34 & -0.26 & 0.60 \\ 0.62 & -0.30 & -0.44 \end{bmatrix}$$

$$E = \begin{bmatrix} -0.33 & -0.24 & 0.57 \\ 0.63 & -0.20 & -0.43 \\ -0.33 & -0.24 & 0.57 \\ 0.60 & -0.21 & -0.39 \end{bmatrix}$$

$Q'(t=0, 1, 2)$ 矩阵为

$$Q^0 = \begin{bmatrix} 1.19 & -1.40 & 0.21 \\ 0.93 & -1.04 & 0.10 \\ 0.70 & -0.79 & 0.09 \\ 1.23 & -1.40 & 0.17 \end{bmatrix}$$

$$Q^1 = \begin{bmatrix} 1.29 & -1.22 & -0.06 \\ 0.98 & -0.94 & -0.04 \\ 0.73 & -0.74 & 0.02 \\ 1.33 & -1.21 & -0.11 \end{bmatrix}$$

$$Q^2 = \begin{bmatrix} 1.33 & -1.06 & -0.27 \\ 1.00 & -0.86 & -0.15 \\ 0.74 & -0.70 & -0.04 \\ 1.36 & -1.04 & -0.32 \end{bmatrix}$$

根据式(17)和式(20),CD对 \hat{X} 和 \hat{A} 的梯度分别为

$$\nabla_{\hat{X}} CD = \begin{bmatrix} -1.24 & 0.34 & -1.14 \\ -1.02 & 0.19 & -0.97 \\ -0.45 & 0.22 & -0.37 \\ -1.44 & 0.27 & -1.37 \end{bmatrix}$$

$$\nabla_{\hat{A}} CD = \begin{bmatrix} -0.30 & -0.18 & -0.65 & -0.26 \\ -0.98 & -0.51 & -0.88 & -0.92 \\ -0.05 & 0.10 & 0.01 & 0.08 \\ -1.50 & -0.81 & -1.53 & -1.43 \end{bmatrix}$$

采用PyTorch的自动求导功能,可得出与上述基于式(17)和式(20)相同的梯度结果.

4.2 梯度攻击算法的可用性实验

为验证本文场景中,邻接矩阵与特征矩阵为离散取值情况下,基于梯度下降法和贪心算法思想设计的扰动算法的可用性,设计一组离散取值的梯度方法实验.实验采用式(13)的Cook距离 $CD(K=2)$.为验证“在一定条件下,基于梯度的算法对离散数据同样有效”这一结论,本组实验将参数矩阵 W^* 视为常数矩阵,以排除图卷积网络重训练环节对梯度有效性的影响.实验中直接求解CD对邻接矩阵 A 的梯度,根据梯度进行扰动,并与遍历方法比较,观察CD值的变化趋势与收敛情况.

实验 1 对数据集中图的连边逐一扰动,并计算Cook距离CD与初始CD之差 ΔCD ,按扰动量从大到小排序;再计算各个扰动连边对应梯度的平均值,将二者绘制于同一图中,如图1所示.再根据算法中的步骤4,由梯度得出使CD值增大的许用梯度,与对应的值的增量,按照 ΔCD 的降序绘制于同一图中,如图2所示.图中所示 ΔCD 与梯度均除以对应的最大值做了归一化处理.本实验基于polblogs数据集完成.

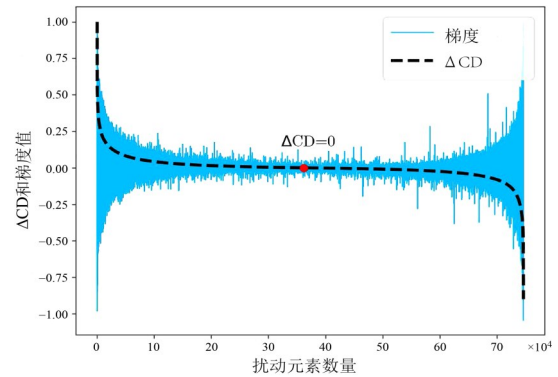


图1 梯度与 ΔCD 的关系图

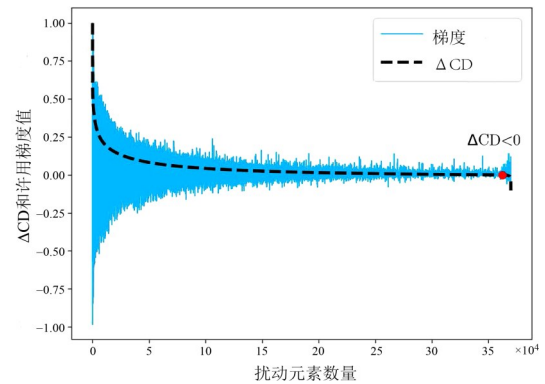


图2 许用梯度与 ΔCD 的关系图

从图1与图2可以看出,扰动一条连边后,CD的实际增幅与梯度之间有强相关关系:对梯度的绝对值大的位置扰动,CD的实际增幅大;对梯度的绝对值小的位置扰动,CD的实际增幅也小.计算图2中扰动量 ΔCD 排序前 n 的 ΔCD 值与对应许用梯度绝对值的Cosine相关系数,如表2所示.可见,CD的实际增量与梯度高度相关,且增幅的排序越靠前(扰动增量越大或梯度绝对值越大),相关性越强.由此得出结论:梯度信息可为对抗攻击场景的离散数据扰动提供重要指导.

表2 ΔCD 与对应许用梯度绝对值的Cosine相关系数

n	10	50	100	200	500
相关系数	0.992	0.992	0.989	0.984	0.971
n	1 000	2000	5 000	10 000	100 000
相关系数	0.971	0.972	0.972	0.969	0.965

由于图2所示的是理论上使CD增长的许用梯度,而红色点后面的数据表示 $\Delta CD < 0$,且图2尾部存在略微的凸起,这部分的数据表明许用梯度标致的扰动方向与CD的实际增减方向不一致.从绝对数量上来看,这些数据占总数的2.016%.事实上, $\Delta CD < -0.05$ 的只有35个,仅占0.0095%,而绝大多数 $\Delta CD < 0$ 即许用梯度产生的错误梯度,为梯度的绝对值极小时,梯度信息已经失去意义而带来的误差.考虑到对抗攻击的实际

场景,实际扰动元素个数不足元素总数的0.1%(连边总数的10%以内),即仅需使用许用梯度绝对值最大的前0.1%的梯度,许用梯度的上述错误对实际攻击场景无任何影响.

实验2 基于梯度下降法与贪心算法的思想,每轮选择许用梯度中梯度的绝对值最大的 γn_0 条连边对邻接矩阵 A 增量式地进行扰动,达到扰动总数 n_0 实验中设置 $n_0=50, 100, 200, 500$, γ 设置为0.04~0.07的随机数).达到扰动总数后,继续对扰动列表进行更新,直至满足迭代总次数iters停止(实验中设置iters=400).记录迭代轮数对应的CD值,并绘制曲线如图3所示.

从实验结果图中可以看出,在4个数据集中,对于

不同的扰动连边数,算法均有收敛的趋势.在一定轮数后,CD值在某一范围内小幅度波动,具有可用性.

从以上2组实验结果来看,实验1表明梯度信息对离散数据扰动有重要理论指导,实验2表明算法有收敛的趋势.实验结果支持了“对于图卷积网络中使用的离散型0,1数据,当确定损失函数后,基于梯度的攻击方法对损失函数的增长具有有效性和可用性”这一结论.由此进一步增加了基于梯度的攻击方法在实用中的说服力.然而究竟在何种条件下能够充分保证这种有效性,以及基于梯度的离散算法在使用时有什么制约条件,需要进一步的理论证明.上述实验结果也对理论研究的可行性提供了证据支持.

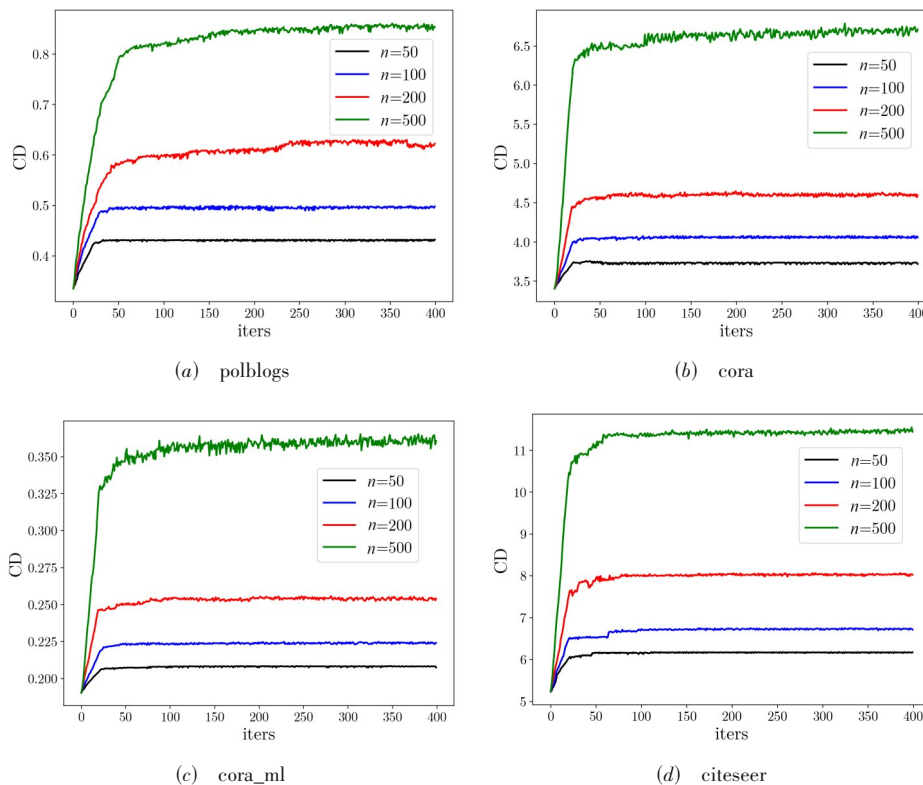


图3 算法1在4个数据集上的收敛性曲线

4.3 参数差异度量的合理性实验

采用式(13)的Cook距离 $CD(K=2, K=3)$ 衡量参数差异.实验中随机生成参数初始值并代入图卷积网络(式(1))训练,迭代次数iters=100.每隔20轮记录1次训练参数,计算CD值,实验结果如表3所示.每组实验CD的最小值用黑体标出.

根据理论分析,随着训练轮数的增加,训练参数应该与最佳参数 W_0 逐渐靠近,即与最佳参数 W_0 的距离应该逐渐减小,从而表现为CD值递减.表中4个数据集在不同K值的实验结果显示,随着训练轮数增加,CD值有减小的趋势,与理论分析一致.

4.4 参数差异假设的合理性实验

本文基于参数差异假设建立攻击模型(式(6)和式(7)),并将现有的其他攻击方法概括为基于矛盾数据假设的攻击方法.为验证参数差异假设的合理性和由其导出方法的有效性,设计以下实验对参数差异假设和矛盾数据假设做对比分析.

实验1 实验中采用本文建立的Cook距离 $CD(K=2, 3)$ 作为基于参数差异假设的攻击方法Parattack.为了与基于矛盾数据假设的攻击方法对比,控制Cook距离表达式的形式不变,基于矛盾数据假设,将表达式

表3 参数差异的度量结果

		0	20	40	60	80	100
polblogs	K=2	348.69	301.92	262.28	228.88	200.81	177.18
	K=3	110.59	95.12	82.02	71.03	61.82	54.10
cora	K=2	635 230.06	3 629.67	1 272.43	695.66	453.81	321.55
	K=3	24 486 698.00	7 124.50	1 825.30	796.38	424.10	251.75
cora_ml	K=2	20 759.96	6 339.43	2 835.57	1 599.68	1 035.22	732.67
	K=3	33 737.08	6 339.66	2 616.48	1 393.28	850.18	567.34
citeseer	K=2	1 063 964.12	4 593.94	1 904.96	1 034.97	607.26	389.56
	K=3	78 890 432.00	2 192 633.00	5 563.29	2 226.66	835.56	385.31

(13)中 $\mathbf{Z} = \mathcal{A}\mathbf{X}$ 更换为 $\hat{\mathbf{Z}} = \hat{\mathcal{A}}\hat{\mathbf{X}}$, 记为 $\widehat{\text{CD}}(K=2,3)$. 实验中划分 40% 训练集和 60% 的测试集, $k=1$ 和 $k=2$ 时, 允许对现有连边总数的 10% 进行攻击. 参数 γ 取 0.03~0.05 的随机数, 算法 1 的迭代次数 $\text{iters}=30$. 重复 10 次实验, 记录 10 次攻击结果准确率下降的平均值, 得到表 4 的实验数据.

从表 4 可以看出, $k=2$ 相较于 $k=1$ 时攻击效果提升明显, 即图卷积网络的攻击难度下降. 控制 Cook 距离 CD 的输入数据为唯一变量, 从整体上该场景中基于参数差异假设的 Cook 距离 CD 的攻击效果好于矛盾数据假设修改的 $\widehat{\text{CD}}$, 部分结果准确率下降幅度成倍数增长. 在部分取值如 $k=1$ 时, cora 和 citeseer 数据集, 基于参数差异假设的攻击效果相比基于矛盾数据假设的攻击效果提升并不显著. 这提示我们: 矛盾数据假设仍是抵抗攻击本质的一种重要概括; 矛盾数据假设是参数差异假设的一种重要实现形式. 在该实验场景下, 基于矛盾数据假设导出的方法, 实质上也是诱导图卷积网络训练出错误参数, 这一观点在文献[22]中同样得到了印证.

表4 参数差异假设与矛盾数据假设导出方法的攻击结果($k=1,2$)

		polblogs	cora_ml	cora	citeseer
$k=1$	CD(K=2)	6.11%	17.33%	21.81%	15.55%
	$\widehat{\text{CD}}(K=2)$	3.21%	9.52%	20.96%	14.31%
	CD(K=3)	5.98%	23.14%	21.34%	15.50%
	$\widehat{\text{CD}}(K=3)$	3.46%	6.59%	20.29%	14.29%
$k=2$	CD(K=2)	14.21%	67.76%	65.30%	51.24%
	$\widehat{\text{CD}}(K=2)$	5.94%	64.51%	64.04%	48.25%
	CD(K=3)	17.26%	70.16%	69.29%	51.00%
	$\widehat{\text{CD}}(K=3)$	5.63%	62.96%	68.24%	46.38%

实验2 设计实验对比分析基于矛盾数据假设的攻击方法 Mettack 和基于参数差异假设的攻击方法 Parattack ($K=2,3$), 以验证参数差异假设的合理性. 在不同扰动量(1%~5%)下, 每种攻击方法均计算自身的损失函数和其他 2 种方法的损失函数, 共 3 种损失函数值, 即计算式(3)的 \mathcal{L} 和式(13)的 $\text{CD}(K=2)$ 和

$\text{CD}(K=3)$. 实验设置与实验 1 相同, 所有数值均为 10 次实验的平均值. 选取 Cora 数据集 $k=1$ 和 $k=2$ 的实验结果如表 5 和表 6 所示. 其他数据集和不同参数下的实验结果与所列结果相似, 实验结论相同, 限于篇幅不再列出.

观察表 5 和表 6 的实验结果. 从横向看, 无论是基于矛盾数据假设的攻击方法 Mettack, 抑或基于参数差异假设的 Parattack ($K=2,3$), 随着扰动比例的增加, 衡量参数差异的 CD 值 ($K=2,3$) 均呈稳定增大趋势. 这表明 Mettack 方法与 Parattack 攻击方法一样, 事实上也是诱导图神经网络训练出异于原始参数的训练参数. 从纵向看, 相同扰动量, Parattack 对应的 3 种损失函数均大于 Mettack 的 3 种损失函数, 说明了基于参数差异假设度量方法的有效性. 值得注意的是, 如表 5 所示, Mettack 方法在 $k=1$ 时损失值 \mathcal{L} 并未稳定增加, 而存在波动. 原因可作以下分析. 损失函数 \mathcal{L} 中有 2 个变量, 即扰动后的邻接矩阵 $\hat{\mathbf{A}}$ 和重训练参数 \mathbf{W}^* . 由 3.1 节可知, 对抗攻击分为两个阶段: 一是根据式(6)实施投毒攻击; 二是根据式(7)进行对抗训练. 投毒攻击阶段基于固定的 \mathbf{W}^* 扰动 $\hat{\mathbf{A}}$ 使 \mathcal{L} 增大; 而对抗训练阶段固定 $\hat{\mathbf{A}}$ 训练 \mathbf{W}^* 使 \mathcal{L} 减小. 最终可能产生两种结果: 若对抗攻击阶段占优, 则最终损失函数 \mathcal{L} 增加; 若对抗训练阶段占优, 则最终损失函数 \mathcal{L} 减小. 两种结果均可出现, 这与表 5 的 Mettack 损失函数 \mathcal{L} 存在波动的实验结果一致. 与损失函数 \mathcal{L} 相比, Cook 距离 CD 中只有一个变量 \mathbf{W}^* , 大幅减弱了该种不一致性而提高攻击损失的灵敏度. 对于 $k=2$ 时, 表 6 的结果显示的损失函数均稳定增加, 对于损失函数 \mathcal{L} 表明其对抗攻击阶段占优, 这与实验 1 得出的 $k=2$ 时图卷积网络的攻击难度下降的结论一致.

以上实验表明, 参数差异假设具有合理性和较为一般的适用性.

4.5 与其他攻击方法的对比实验

本节将本文提出的基于参数差异假设的攻击方法 Parattack ($K=2,3$) 与基准方法进行对比实验. 基准方法包括随机攻击方法 Random 和基于矛盾数据假设的攻

表5 参数差异假设的合理性验证(Cora数据集,SGC模型取 $k=1$)

方法		未扰动	1%	2%	3%	4%	5%
Mettack	\mathcal{L}	44.19	74.09	67.75	59.41	53.27	46.71
	CD ($K=2$)	9.81	303.06	794.30	1 571.02	3 135.19	5 625.58
	CD ($K=3$)	2.28	3 865.35	15 574.70	37 658.80	89 857.40	178 598.00
Parattack ($K=2$)	\mathcal{L}	44.01	90.22	98.59	103.03	102.39	99.31
	CD ($K=2$)	9.75	402.82	1 754.35	4 355.42	6 491.25	8 567.13
	CD ($K=3$)	2.26	11 776.00	93 541.30	346 686.00	463 502.00	583 897.00
Parattack ($K=3$)	\mathcal{L}	44.19	88.99	95.03	100.08	109.78	100.64
	CD ($K=2$)	9.80	312.38	1 314.19	3 230.10	8 447.78	7 479.64
	CD ($K=3$)	2.27	9 956.53	69 866.20	218 892.00	745 062.00	572 026.00

表6 参数差异假设的合理性验证(Cora数据集,SGC模型取 $k=2$)

方法		未扰动	1%	2%	3%	4%	5%
Mettack	\mathcal{L}	88.50	132.89	1 090.86	9 216.20	54 523.90	144 609.00
	CD ($K=2$)	13.97	243.53	3 854.18	13 369.40	27 487.30	43 982.20
	CD ($K=3$)	2.59	5 265.44	103 406.00	256 005.00	306 480.00	441 769.00
Parattack ($K=2$)	\mathcal{L}	88.08	212.14	4 489.40	36 598.90	138 945.00	381 598.00
	CD ($K=2$)	13.86	675.67	9 962.60	26 537.40	73 422.50	112 915.00
	CD ($K=3$)	2.57	26 488.00	208 945.00	439 317.00	940 865.00	1 004 410.00
Parattack ($K=3$)	\mathcal{L}	88.29	200.95	5 794.97	44 233.90	147 960.00	425 041.00
	CD ($K=2$)	13.92	1 175.71	11 471.30	36 288.10	50 882.00	141 631.00
	CD ($K=3$)	2.58	55 506.60	244 881.00	645 585.00	536 679.00	1 239 370.00

击方法 DICE, Mettack. 下面分别对基准方法简要介绍.

(1) 随机攻击方法 Random. 该方法从训练集中随机选择连边增加或删除.

(2) “删除同类连接异类”攻击方法 DICE (Delete Internally Connect Externally). 该方法根据“从同类节点中删除连边, 在不同类节点间增加连边”这一规则, 随机增删连边. 根据本文观点, 该方法属于基于矛盾数据假设的攻击方法.

(3) 基于元学习的攻击方法 Mettack. 该攻击方法的损失函数与训练图卷积网的损失函数相同, 损失函数代入扰动后的投毒训练数据, 可由矛盾数据假设导出; 然后根据损失函数得攻击梯度, 并采用贪心算法进行扰动筛选.

实验将数据集随机地划分为 40% 的训练集和 60% 的测试集, 允许使用 3% 的连边数进行攻击. 根据攻击梯度, Mettack 方法与 Parattack 方法均使用离散梯度下降算法, 参数 γ 取 0.03~0.05 的随机数, 算法的迭代次数 $iters=50$. 对邻接矩阵中的连边进行扰动, 重复 10 次实验, 记录 10 次攻击结果准确率下降的平均值, 得到表 7 的实验数据. 同时表 8 列出了不同 k 值时, 图卷积网络 SGC 在各个未受到对抗攻击污染的数据集上的准确率. 实验结果表明, 随机攻击方法 Random 和基于规则的 DICE 几乎无法实施有效的数据投毒攻击, 即对重训练的 SGC 几乎没有攻击效果; 部分实验结果为负值, 攻

表7 Parattack(CD)与其他方法的攻击结果

		polblogs	cora_ml	cora	citeseer
$k=1$	Random	0.07%	-0.16%	0.49%	-0.18%
	DICE	0.06%	0.07%	0.09%	-0.13%
	Mettack	3.99%	2.46%	12.85%	4.83%
	CD($K=2$)	4.36%	6.26%	15.62%	6.57%
	CD($K=3$)	3.80%	6.79%	11.13%	5.35%
$k=2$	Random	-0.11%	0.03%	0.13%	4.07%
	DICE	0.25%	0.05%	-0.45%	0.06%
	Mettack	3.67%	40.75%	46.00%	23.31%
	CD($K=2$)	3.12%	52.50%	46.48%	29.26%
$k=3$	CD($K=3$)	4.40%	48.69%	43.04%	29.19%
	Random	-0.15%	-0.02%	0.15%	0.33%
	DICE	0.31%	0.03%	0.01%	2.77%
	Mettack	30.83%	53.96%	58.66%	27.22%
	CD($K=2$)	42.16%	56.99%	63.59%	52.82%
CD($K=3$)	30.07%	56.74%	64.29%	49.93%	

表8 图卷积网络 SGC 在未受攻击的数据集上的准确率

	polblogs	cora_ml	cora	citeseer
$k=1$	91.84%	86.69%	85.20%	76.44%
$k=2$	94.88%	83.36%	86.13%	77.82%
$k=3$	94.79%	82.77%	86.19%	77.56%

击后准确率不减反增, 但多为随机波动. 原因是该两种方法或随机选取连边, 或基于规则随机增删连边, 未进

重训练. 缺少针对性的训练环节是导致攻击无效的主要原因. 对于具备重训练环节的 Mettack 和 Parattack, 实验结果表明, 基于参数差异假设的攻击方法 Parattack 相比基于矛盾数据假设的攻击方法 Mettack 更有效性. 大部分攻击效果都有较大幅度提升, 从而验证了基于参数差异假设导出方法的有效性. 仅攻击特征以及结构和特征联合的攻击方法, 与上述方法完全一致, 实验结果大致相仿, 不再赘述.

5 结论

本文基于参数差异假设, 将统计诊断的重要研究成果 Cook 距离引入图卷积网络对抗性攻击, 用以衡量模型的参数差异. 本文主要得出了以下结论:

一是对于邻接矩阵为 0, 1 的离散数据场景, 基于梯度的攻击方法具有有效性和可用性. 基于攻击梯度的离散数据实验结果支持了这一结论. 实验结果也为算法的收敛性等理论研究的可行性提供了证据支持.

二是推导得出 Cook 距离的计算式和基于 Cook 距离的攻击梯度闭式解. 攻击梯度的闭式解为进一步寻找攻击有效性的本质和简化攻击方法提供了必要前提.

三是设计了矛盾数据假设和参数差异假设的对比实验, 说明参数差异假设的合理性. 基于参数差异假设的攻击方法与其他方法的对比实验结果说明了参数差异假设导出方法的有效性.

图卷积网络乃至图神经网络对抗性攻击尤其是对抗性攻击存在性本质, 仍具有丰富而深刻的问题有待进一步研究. 例如, 参数差异假设是否包含矛盾数据假设, 矛盾数据假设是否必然蕴含参数差异假设, 是否存在仅满足矛盾数据假设的对抗性攻击而不满足参数差异假设, 或是否存在满足参数差异假设的数据样本而不属于对抗性攻击样本, 矛盾数据假设与参数差异假设何时等价, 对抗性攻击问题本质能否更紧凑地定量表达, 等等, 需要进一步的理论和实验研究.

附录

定理 1 的证明.

$$d(\nabla_W \mathcal{L}(W)) = Z^T d[\sigma(ZW)] I_m$$

$$\begin{aligned} \text{vec}[d(\nabla_W \mathcal{L}(W))] &= (I_m \otimes Z^T) \text{vec}[d\sigma(ZW)] \\ &= (I_m \otimes Z^T) \frac{\partial \sigma^T(ZW)}{\partial (ZW)} (I_m \otimes Z) \text{vec}(dW) \end{aligned}$$

代入 Cook 距离的定义式(12)可得定理 1.

证毕

定理 2 的证明.

设 $v(W^*) = \text{vec}[\nabla_W \mathcal{L}(W^*)]$, 则

$$\text{vec}(\nabla_{\hat{X}} \text{CD}) = \frac{\partial \hat{Z}}{\partial \hat{X}} \frac{\partial W^*}{\partial \hat{Z}} \frac{\partial v(W^*)}{\partial W^*} \frac{\partial \text{CD}}{\partial v(W^*)}$$

其中, $\frac{\partial \hat{Z}}{\partial \hat{X}}, \frac{\partial W^*}{\partial \hat{Z}}, \frac{\partial v(W^*)}{\partial W^*}, \frac{\partial \text{CD}}{\partial v(W^*)}, \nabla_{\hat{X}} \text{CD}$ 按如下定

义, 设矩阵 F , X , 有

$$\frac{\partial F}{\partial X} = \frac{\partial \text{vec}(F)}{\partial \text{vec}(X)}, \quad \nabla_X \text{CD} = \left[\frac{\partial \text{CD}}{\partial X_{ij}} \right]$$

(1) 求 $\frac{\partial \text{CD}}{\partial v(W^*)}$

$$\begin{aligned} \frac{\partial \text{CD}}{\partial v(W^*)} &= 2(\nabla_W^2 \mathcal{L}(W^*))^{K-2} v(W^*) \\ &= -2M^{K-2}(W^*) \text{vec}(Z^T E) \end{aligned}$$

(2) 求 $\frac{\partial v(W^*)}{\partial W^*}$

$$\frac{\partial v(W^*)}{\partial W^*} = (I_m \otimes Z^T) K_{mN} \Lambda^* K_{Nm} (I_m \otimes Z) = M(W^*)$$

$$\frac{\partial v(W^*)}{\partial W^*} \frac{\partial \text{CD}}{\partial v(W^*)} = -2M^{K-1}(W^*) \text{vec}(H_1) = -2\text{vec}(H_K)$$

(3) 求 $\frac{\partial W^*}{\partial \hat{Z}}$

$$\begin{aligned} d\text{vec}\left\{\hat{Z}^T[\sigma(\hat{Z}W^t) - Y]\right\} \\ &= d\left\{(I_m \otimes \hat{Z}^T) \text{vec}\left[\sigma(\hat{Z}W^t) - Y\right]\right\} \\ &= dv_1(\hat{Z}) + dv_2(\hat{Z}) \end{aligned}$$

$$dv_1(\hat{Z}) = K_{mn} \left\{ I_n \otimes [\sigma(\hat{Z}W^t) - Y]^T \right\} \text{vec}(d\hat{Z})$$

$$dv_2(\hat{Z}) = (I_m \otimes \hat{Z}^T) K_{mN} \hat{\Lambda} K_{Nm} [(W^t)^T \otimes I_N] \text{vec}(d\hat{Z})$$

$$\frac{\partial \text{vec}\left\{\hat{Z}^T[\sigma(\hat{Z}W^t) - Y]\right\}}{\partial \hat{Z}}$$

$$\begin{aligned} &= \left[I_n \otimes (\sigma(\hat{Z}W^t) - Y) \right] K_{nm} \\ &\quad + (W^t \otimes I_N) K_{mN} \hat{\Lambda} K_{Nm} (I_m \otimes \hat{Z}) \end{aligned}$$

$$\frac{\partial W^*}{\partial \hat{Z}} = \frac{-\alpha \sum_{t=0}^{t_0-1} \partial \hat{Z}^T [\sigma(\hat{Z}W^t) - Y]}{\partial \hat{Z}}$$

$$\begin{aligned} &= -\alpha \sum_{t=0}^{t_0-1} \left[I_n \otimes (\sigma(\hat{Z}W^t) - Y) \right] K_{nm} \\ &\quad + (W^t \otimes I_N) K_{mN} \hat{\Lambda} K_{Nm} (I_m \otimes \hat{Z}) \end{aligned}$$

$$\begin{aligned} & \frac{\partial \mathbf{W}^*}{\partial \hat{\mathbf{Z}}} \frac{\partial v(\mathbf{W}^*)}{\partial \mathbf{W}^*} \frac{\partial \text{CD}}{\partial v(\mathbf{W}^*)} \\ &= 2\alpha \left\{ \sum_{t=0}^{t_0-1} \left[\mathbf{I}_n \otimes (\sigma(\hat{\mathbf{Z}} \mathbf{W}^t) - \mathbf{Y}) \right] \mathbf{K}_{nm} \right. \\ & \quad \left. + (\mathbf{W}^t \otimes \mathbf{I}_N) \mathbf{K}_{mN} \hat{\mathbf{A}}^t \mathbf{K}_{Nm} (\mathbf{I}_m \otimes \hat{\mathbf{Z}}) \right\} \text{vec}(\mathbf{H}_K) \\ &= 2\alpha \text{vec} \left(\sum_{t=0}^{t_0-1} \mathbf{Q}^t \mathbf{W}^{t\top} - \hat{\mathbf{E}}^t \mathbf{H}_K^{\top} \right) \end{aligned}$$

其中, $\text{vec}(\mathbf{Q}^t) = \hat{\mathbf{A}}^t \text{vec}(\mathbf{H}_K^{\top} \hat{\mathbf{Z}}^t)$.

(4) 求 $\frac{\partial \hat{\mathbf{Z}}}{\partial \hat{\mathbf{X}}}$

$$\begin{aligned} \frac{\partial \hat{\mathbf{Z}}}{\partial \hat{\mathbf{X}}} &= \mathbf{I}_n \otimes \mathcal{A} \\ \text{vec}(\nabla_{\hat{\mathbf{X}}} \text{CD}) &= \frac{\partial \hat{\mathbf{Z}}}{\partial \hat{\mathbf{X}}} \frac{\partial \mathbf{W}^*}{\partial \hat{\mathbf{Z}}} \frac{\partial v(\mathbf{W}^*)}{\partial \mathbf{W}^*} \frac{\partial \text{CD}}{\partial v(\mathbf{W}^*)} \\ &= 2\alpha (\mathbf{I}_n \otimes \hat{\mathbf{A}}) \text{vec} \left(\sum_{t=0}^{t_0-1} \mathbf{Q}^t \mathbf{W}^{t\top} - \hat{\mathbf{E}}^t \mathbf{H}_K^{\top} \right) \\ \nabla_{\hat{\mathbf{X}}} \text{CD} &= 2\alpha \cdot \hat{\mathbf{A}} \sum_{t=0}^{t_0-1} (\mathbf{Q}^t \mathbf{W}^{t\top} - \hat{\mathbf{E}}^t \mathbf{H}_K^{\top}) \end{aligned}$$

证毕

定理3的证明.

用归纳法容易证明,有

$$\begin{aligned} \frac{\partial \hat{\mathbf{Z}}}{\partial \hat{\mathbf{A}}} &= \left(\hat{\mathbf{D}}^{-\frac{1}{2}} \otimes \hat{\mathbf{D}}^{-\frac{1}{2}} \right) \left[\sum_{i=0}^{k-1} (\hat{\mathbf{L}}^{k-i-1} \hat{\mathbf{X}}) \otimes \hat{\mathbf{L}}^i \right] \\ \text{vec}(\nabla_{\hat{\mathbf{A}}} \text{CD}) &= \frac{\partial \hat{\mathbf{Z}}}{\partial \hat{\mathbf{A}}} \frac{\partial \mathbf{W}^*}{\partial \hat{\mathbf{Z}}} \frac{\partial v(\mathbf{W}^*)}{\partial \mathbf{W}^*} \frac{\partial \text{CD}}{\partial v(\mathbf{W}^*)} \\ &= 2\alpha \frac{\partial \hat{\mathbf{Z}}}{\partial \hat{\mathbf{A}}} \text{vec} \left(\sum_{t=0}^{t_0-1} \mathbf{Q}^t \mathbf{W}^{t\top} - \hat{\mathbf{E}}^t \mathbf{H}_K^{\top} \right) \\ &= 2\alpha \text{vec} \left\{ \hat{\mathbf{D}}^{-\frac{1}{2}} \left[\sum_{i=0}^{k-1} \hat{\mathbf{L}}^i \left(\sum_{t=0}^{t_0-1} \mathbf{Q}^t \mathbf{W}^{t\top} - \hat{\mathbf{E}}^t \mathbf{H}_K^{\top} \right) \right. \right. \\ & \quad \left. \left. (\hat{\mathbf{L}}^{k-i-1} \hat{\mathbf{X}})^{\top} \right] \hat{\mathbf{D}}^{-\frac{1}{2}} \right\} \\ \nabla_{\hat{\mathbf{A}}} \text{CD} &= 2\alpha \hat{\mathbf{D}}^{-\frac{1}{2}} \left[\sum_{i=0}^{k-1} \hat{\mathbf{L}}^i \left(\sum_{t=0}^{t_0-1} \mathbf{Q}^t \mathbf{W}^{t\top} - \hat{\mathbf{E}}^t \mathbf{H}_K^{\top} \right) \right. \\ & \quad \left. (\hat{\mathbf{L}}^{k-i-1} \hat{\mathbf{X}})^{\top} \right] \hat{\mathbf{D}}^{-\frac{1}{2}} \end{aligned}$$

证毕

参考文献

[1] SZEGEDY C, ZAREMBA W, SUTSKEVER I, et al. Intriguing properties of neural networks[EB/OL]. (2013)

[2021]. <https://arxiv.org/abs/1312>.

- [2] KURAKIN A, GOODFELLOW I, BENGIO S. Adversarial examples in the physical world[EB/OL]. (2016)[2021]. <https://arxiv.org/abs/1607.02533v1>.
- [3] ILYAS A, SANTURKAR S, TSIPRAS D, et al. Adversarial examples are not bugs, they are features[C]//Advances in Neural Information Processing Systems. Vancouver: NIPS, 2019: 125-136.
- [4] YUAN X Y, HE P, ZHU Q L, et al. Adversarial examples: Attacks and defenses for deep learning[J]. IEEE Transactions on Neural Networks and Learning Systems, 2019, 30(9): 2805-2824.
- [5] COSTA L DA F, RODRIGUES F A, TRAVIESO G, et al. Characterization of complex networks: A survey of measurements[J]. Advances in Physics, 2007, 56(1): 167-242.
- [6] KIPF T N, WELING M. Semi-supervised classification with graph convolutional networks[EB/OL]. (2017)[2021]. <https://arxiv.org/abs/1609.02907v4>.
- [7] WU F, ZHANG T, SOUZA JR A H d, et al. Simplifying graph convolutional networks[EB/OL]. (2019) [2021]. <https://arxiv.org/abs/1902.07153v1>.
- [8] LI Q M, WU X M, LIU H, et al. Label efficient semi-supervised learning via graph filtering[C]//2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). Long Beach: IEEE, 2019: 9574-9583.
- [9] NT H, MAEHARA T. Revisiting graph neural networks: All we have is low-pass filters[EB/OL]. (2019) [2021]. <https://arxiv.org/abs/1905.09550v1>.
- [10] SCARSELLI F, GORI M, TSOI A C, et al. The graph neural network model[J]. IEEE Transactions on Neural Networks, 2009, 20(1): 61-80.
- [11] WU Z H, PAN S R, CHEN F W, et al. A comprehensive survey on graph neural networks[J]. IEEE Transactions on Neural Networks and Learning Systems, 2021, 32(1): 4-24.
- [12] 徐冰冰, 岑科廷, 黄俊杰, 等. 图卷积神经网络综述[J]. 计算机学报, 2020, 43(5): 755-780.
- XU B B, CEN K T, HUANG J J, et al. A survey on graph convolutional neural network[J]. Chinese Journal of Computers, 2020, 43(5): 755-780. (in Chinese)
- [13] 白铂, 刘玉婷, 马驰骋, 等. 图神经网络[J]. 中国科学: 数学, 2020, 50(3): 367-384.
- BAI B, LIU Y T, MA C C, et al. Graph neural network [J]. Scientia Sinica(Mathematica), 2020, 50(3): 367-384. (in Chinese)
- [14] ZÜGNER D, AKBARNEJAD A, GÜNNEMANN S. Ad-

- versarial attacks on neural networks for graph data[C]//KDD' 18: Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining. London: ACM, 2018: 2847-2856.
- [15] ENTEZARI N, AL-SAYOURI S A, DARVISHZADEH A, et al. All You need is low (rank): Defending against adversarial attacks on graphs[C]//WSDM'20: Proceedings of the 13th International Conference on Web Search and Data Mining. Houston: ACM, 2020: 169-177.
- [16] LI J, ZHANG H L, HAN Z C, et al. Adversarial attack on community detection by hiding individuals[C]//WWW' 20: Proceedings of The Web Conference 2020. Taipei: ACM, 2020: 917-927.
- [17] WU Y T, LIU W, HU X B, et al. Parameter discrepancy hypothesis: Adversarial attack for graph data[J]. Information Sciences, 2021, 577: 234-244.
- [18] CHEN L, LI J T, PENG JY, et al. A survey of adversarial learning on graphs[EB/OL]. (2020) [2021]. <https://arxiv.org/abs/2003.05730>.
- [19] JIN W, LI Y, XU H, et al. Adversarial attacks and defenses on graphs: A review and empirical study[EB/OL]. (2020)[2021]. <https://arxiv.org/abs/2003.00653v2>.
- [20] LI Y, JIN W, XU H, et al. DeepRobust: A PyTorch library for adversarial attacks and defenses[EB/OL]. (2020) [2021]. <https://arxiv.org/abs/2005.06149>.
- [21] BOJCHEVSKI A, GÜNNEMANN S. Adversarial attacks on node embeddings via graph poisoning[C]//International Conference on Machine Learning. Long Beach: PMLR, 2019: 695-704.
- [22] ZÜGNER D, GÜNNEMANN S. Adversarial attacks on graph neural networks via meta learning[C]//International Conference on Learning Representations. New Orleans: ICLR, 2019: 1112.
- [23] COOK R D. Detection of influential observation in linear regression[J]. Technometrics, 1977, 19(1): 15-18.
- [24] COOK R D. Influential observations in linear regression [J]. Journal of the American Statistical Association, 1979, 74(365): 169-174.
- [25] COOK R D, WEISBERG S. Residuals and Influence in Regression[M]. New York: Chapman and Hall, 1982.
- [26] WEI B C, SHIH J Q. On statistical models for regression diagnostics[J]. Annals of the Institute of Statistical Mathematics, 1994, 46(2): 267-278.
- [27] 费宇, 陈飞, 喻达磊. 线性和广义线性混合模型及其统计诊断[M]. 北京: 科学出版社, 2013.
- [28] 韦博成, 林金官, 解锋昌. 统计诊断[M]. 北京: 高等教育

出版社, 2009.

WEI B C, LIN J G, XIE F C. Statistical Diagnostics[M]. Beijing: Higher Education Press, 2009. (in Chinese)

- [29] 韦博成, 鲁国斌, 史建清. 统计诊断引论[M]. 南京: 东南大学出版社, 1991.
- [30] SEN P, NAMATA G, BILGIC M, et al. Collective classification in network data[J]. AI Magazine, 2008, 29(3): 93.
- [31] MCCALLUM A K, NIGAM K, RENNIE J, et al. Automating the construction of Internet portals with machine learning[J]. Information Retrieval, 2000, 3(2): 127-163.
- [32] ADAMIC L A, GLANCE N. The political blogosphere and the 2004 US election: divided they blog[C]//Proceedings of the 3rd International Workshop on Link Discovery. New York: ACM, 2005: 36-43.

作者简介



吴翼腾 男, 1992年出生, 吉林省吉林人. 博士, 工程师. 主要研究方向为网络空间安全、人工智能安全、对抗机器学习.
E-mail: wuyiteng1992@163.com



刘伟 男, 1992年出生, 河北人保定人. 硕士, 工程师. 主要研究方向为人工智能安全、自然语言处理.



于淑乔 女, 1998年生, 河南郑州人. 硕士. 主要研究方向为数据科学.