

# 一种防御投影梯度下降攻击的图时空注意力网络方法

尹艺, 吴杨\*, 张权, 曾卓, 李娅洁

重庆电子科技职业大学人工智能与大数据学院, 重庆

收稿日期: 2026年4月23日; 录用日期: 2026年5月22日; 发布日期: 2026年5月29日

## 摘要

图时空网络在交通预测与疫情分析等领域广泛应用, 然而其时空耦合结构也使其易受对抗攻击干扰。针对投影梯度下降PGD对抗攻击对模型性能的破坏, 文章提出了一种鲁棒图时空自注意力网络RGSTAN, 该方法引入交叉时空自注意力强化时空特征建模能力, 并结合时空平滑策略在空间与时间维度上对动态节点特征表示进行平滑处理与对抗扰动, 降低对抗噪声对关键特征的影响, 从而提升模型鲁棒性。此外, 还分析了图时空模型面临的安全威胁, 揭示了投影梯度下降对抗攻击在时空信息传播过程中的干扰原理。最后, 基于疫情与流量动态图数据的实验结果表明, 在不同攻击约束的对抗环境下, 时空平滑机制能强化时空自注意力网络层的鲁棒性, 使RGSTAN能有效缓解投影梯度攻击的干扰。

## 关键词

图时空网络, 投影梯度下降, 时空自注意力, 时空平滑

# A Graph Spatial-Temporal Attention Network Method for Defending Against Projected Gradient Descent Attacks

Yi Yin, Yang Wu\*, Quan Zhang, Zhuo Zeng, Yajie Li

School of Artificial Intelligence and Big Data, Chongqing Polytechnic University of Electronic Technology, Chongqing

Received: April 23, 2026; accepted: May 22, 2026; published: May 29, 2026

\*通讯作者。

文章引用: 尹艺, 吴杨, 张权, 曾卓, 李娅洁. 一种防御投影梯度下降攻击的图时空注意力网络方法[J]. 计算机科学与应用, 2026, 16(5): 493-507. DOI: 10.12677/csa.2026.165200

## Abstract

Graph spatial-temporal networks are widely used for traffic forecasting and epidemic analysis, but their spatial-temporal correlated structure makes them susceptible to perturbation from adversarial attacks. In order to limit the performance degradation caused by projected gradient descent attacks, this paper proposes a Robust Graph Spatial-Temporal Self-Attention Network model. The model integrates cross-spatial-temporal self-attentions to enhance the spatial-temporal representation dependencies, and utilizes a spatial-temporal smoothing strategy that smooths and perturbs dynamic node features from spatial and temporal dimensions, thus mitigating the impact of adversarial noise on crucial representations and improving the model's robustness. Additionally, this paper explores the security threats of graph spatial-temporal models and analyzes the PGD attack mechanism, which interferes with spatial-temporal information flow. Finally, experimental results based on epidemic and traffic dynamic graph data show that the spatiotemporal smoothing mechanism can enhance the robustness of the spatiotemporal self-attention network layer under adversarial environments with different attack constraints, enabling RGSTAN to effectively mitigate the interference of projection gradient attacks.

## Keywords

Graph Spatial-Temporal Network, Projected Gradient Descent, Spatial-Temporal Self-Attention, Spatial-Temporal Smoothing

Copyright © 2026 by author(s) and Hans Publishers Inc.

This work is licensed under the Creative Commons Attribution International License (CC BY 4.0).

<http://creativecommons.org/licenses/by/4.0/>



Open Access

## 1. 引言

当前社会的各个领域产生了大量随时间演化的关系型数据,例如人际关系、交通网络等[1]。图神经网络作为提取关系型数据的结构化信息的主要方法,能够与循环神经网络 RNN、长短记忆 LSTM 等方法相结合为图时空网络(Graph Spatial-Temporal Network, GSTN) [2]-[4],从随时间演化的图数据中提取由图结构和图演化信息结合的时空特征,常用于执行社交关系推理、交通流量预测等任务。然而,图时空网络的时空耦合特性反而容易引起安全问题,错误的空间信息聚合会干扰邻居节点信息,而错误的时间模式干扰未来的预测结果,故而攻击者只需改动部分时刻的少量节点,就可能严重影响模型性能。

模型对抗攻击能够在模型的输入上添加扰动使模型输出发生错误。其中,投影梯度下降 PGD (Projected Gradient Descent) [5]对抗攻击可以在输入空间中查找使模型输出错误的最小扰动,它能够通过图时空网络的时空信息的传递能力将部分节点的错误从时空方向上传播,攻击者仅需要在节点特征上施加微小扰动,即可通过时空耦合特性将误差逐级放大。本研究提出了鲁棒图时空自注意力网络(Robust Graph Spatial-Temporal Attention Network, RGSTAN),实现了图时空模型对 PGD 攻击的鲁棒优化,并将其与三种防御策略(对抗训练、防御蒸馏、输入净化)下的各类时空模型进行对比,以证明其效用。

本文主要贡献如下:

- 1) 本文构建了面向图时空特征的交叉自注意力融合模型结构 RGSTAN,该模型包含了强化空间

层、强化时间层、时空融合层以及线性输出层。空间层基于图注意力网络实现邻居信息聚集，时间层采用长短注意力方式适应性地汇聚图演化信息，时空融合层采用交叉自注意力机制实现了时空信息的动态融合。

2) 本文提出了时空平滑机制 STS (Spatial-Temporal Smoothing)，其中，空间平滑机制调控节点特征与领域信息之间的权衡以缓解局部特征偏移，时间平滑机制减少对抗扰动引起的图演化特征异常，再通过模型输出添加对抗性噪声，从而提升 RGSTAN 模型的鲁棒性。

3) 本文基于 Chickenpox、EngCovid 与 Wiki-Math 三份公开动态图数据集，将所提出的 RGSTAN 模型与采用不同防御方式的四种主流图时空网络进行比较，在不同扰动约束范围的 PDG 攻击上开展模型的鲁棒优化验证。

## 2. 相关工作

### 2.1. 图时空网络安全威胁分析

图时空网络面临对抗攻击导致的模型精度下降问题。先兴平[6]指出，攻击者可以通过恶意更改图数据来影响模型的性能，包括通过增删边、篡改特征实现投毒，在训练阶段搜索最优结构扰动，以及直接向隐藏状态注入梯度噪声等，从而导致金融风险、商务欺诈、社交欺骗等严重后果。金柯君[7]以投影梯度下降为基础，在模型训练阶段扰动图结构，使模型在学习攻击中受到持续干扰。柏杨[8]在图垂直联邦学习中采用了以节点特征为目标的攻击策略，随机加入噪声以干扰模型精度。目前，这些方法都着重于通过多轮或迭代式的方式干扰模型。故而，本文研究投影梯度攻击对图时空模型的安全影响，并针对性地提升模型对迭代式对抗攻击的鲁棒性。

### 2.2. 图时空网络鲁棒优化研究

目前，针对 PGD 攻击的防御方法主要围绕模型训练或者数据处理方面来提升模型鲁棒性。对抗训练 AT (Adversarial Learning) [9]可由 PGD 扰动的对抗样本与原样本混合训练提升模型对于时空信息扰动的鲁棒性，该方法在训练阶段实时扰动训练样本，容易过度干扰模型的精度，防御蒸馏 DD (Defensive Distillation) [10]通过软标签的方式让模型的决策边界更平滑，降低模型对输入扰动的敏感性。输入净化 IP (Input Purification) [11]在模型训练前先对含对抗扰动的数据样本进行净化，将其重构为更接近原始时空信号的干净输入，再用净化后的样本进行模型训练以降低其扰动影响。王煜恒[12]提出面向图数据的随机平滑机制，在原图结构上引入随机扰动然后进行平滑操作，能够在一定扰动范围内保持预测精度。然而，这些方法多基于单个图快照或者图局部的扰动，缺乏对于跨时间累积扰动效应的建模，难以适用于面向图时空模型的 PGD 攻击。PGD 造成的扰动能够在时空耦合维度上进行优化，并根据动态图数据的时空依赖逐步累积扰动效应。

## 3. 面向动态图数据的图时空网络模型

### 3.1. 动态图数据

随时间演化的关系型数据可以表示为动态图数据  $G^{1 \leq t \leq T} = (V^{1 \leq t \leq T}, E^{1 \leq t \leq T})$ ，它由时间连续的  $T$  个关系型图快照所组成。其中，图快照  $G^t$  包含节点集合  $V^t = \{v'_1, v'_2, \dots, v'_n\}$  和边集合  $E^t = \{e'_1, e'_2, \dots, e'_m\}$ ，以及节点  $c$  维度标签分布  $Y^t = \{y'_i \mid v'_i \in V^t\} \in \mathbb{R}^{n \times c}$ 。

在图快照中，节点特征与节点间的关联矩阵可以表示为  $(X^{1 \leq t \leq T} \in \mathbb{R}^{n \times d}, A^{1 \leq t \leq T} \in \mathbb{R}^{n \times n})$ ，其中， $n$  为节点数量， $d$  为节点特征的维度。

### 3.2. 图时空网络模型

图时空网络模型将具有图结构信息获取能力的 GNN 与具有时间演化信息提取能力的 RNN 相结合，其具备了从动态图数据获取提取时空特征并生成动态图嵌入的能力，并能够基于所提取的动态图嵌入来执行节点分类、边预测等任务。

#### 3.2.1. 空间网络层

图时空网络模型的空间层可通过 GNN 从节点中学习节点特征，并将邻居节点信息聚合并更新节点特征。最具代表性的空间网络层模型为图卷积神经网络 GCN，该网络通过拉普拉斯矩阵从图数据中获取多跳邻居的信息以形成节点的特征表示。

$$H_i^{(l+1)} = \sigma \left( \tilde{D}^{-\frac{1}{2}} \tilde{A} \tilde{D}^{-\frac{1}{2}} H_i^{(l)} W^{(l)} \right) \quad (1)$$

其中  $\sigma$  是激活函数， $H_i^{(0)} = X^t$ ， $\tilde{A} = A + I$ ， $\tilde{D}_{ii} = \sum_j \tilde{A}_{ij}$ ， $W^{(l)} \in \mathbb{R}^{d \times d}$ 。

#### 3.2.2. 时间网络层

动态图中图结构信息的演化不仅依赖当前节点状态，还依赖之前的节点状态。时间网络层能够在时间序列实现共享参数，并将当前快照的图结构信息传递到下一个快照的图结构信息上从而生成具备图演化信息的图嵌入。例如，循环神经网络 RNN 能够通过隐藏状态来实现图结构信息在时间上的传递：

$$Z_t = f \left( W_h \left[ Z_{t-1}, H_t^{(L)} \right] + b \right) \quad (2)$$

其中， $Z_{t-1} \in \mathbb{R}^{n \times d_z}$  表示当前快照  $t-1$  中节点的隐藏状态且维度为  $(n, d_z)$ ， $H_t^{(L)}$  为最后一层空间网络层的输出， $[\ ]$  将两个向量可以拼接为一个新的向量，并通过时间层的权重  $W_h \in \mathbb{R}^{d_z \times (d_z + d)}$  与偏置向量  $b \in \mathbb{R}^{d_z}$  实现特征转化。

### 3.3. 模型损失函数

在时间层输出节点表示  $Z_t$  后，图时空网络模型的最后一层线性层能够将  $Z_t$  转化为预测结果  $\hat{Y}^t$ ，并计算模型的损失函数以应对不同的任务：

1) 节点分类任务的损失函数  $L_{cl}$ ：

$$\hat{Y}^t = \text{softmax}(W_o Z_t + b_o); L_{cl} = - \sum_{t=1}^T \sum_{i=1}^n \sum_{c=1}^C y_{ic}^t \log(p_{ic}^t) \quad (3)$$

其中， $W_o$  为线性层的权重向量， $\hat{Y}^t$  为线性层的输出， $b_o$  为其偏置向量， $y_{ic}^t$  表示节点  $i$  是否属于类别  $c$ ， $p_{ic}^t$  表示模型预测节点  $i$  属于类别  $c$  的概率。

2) 节点回归任务的损失函数  $L_{re}$ ：

$$\hat{Y}^t = W_o Z_t + b_o; L_{re} = - \frac{1}{T \times n} \sum_{t=1}^T \sum_{i=1}^n (y_i^t - \hat{y}_i^t) \quad (4)$$

在节点回归任务重，损失函数采用均方误差来计算真实值与预测值之间的平方误差，再在时间与节点上取均值，通过最小化误差可以使图时空模型学习到更精确预测结果。

## 4. 鲁棒的图时空网络模型

针对投影梯度下降攻击带来的模型性能干扰，本文基于自注意力机制的动态加权聚合[13]，提出了鲁棒图时空自注意力网络(Robust Graph Spatial-Temporal Attention Network, RGSTAN)，该模型(见图 1)不但

捕获空间维度上的图结构关系以及时间维度上的动态演化关系，还能通过引入时空特征平滑以缓解对时空特征的破坏。该模型主要由强化空间层、强化时间层、时空融合层与线性输出层所组成，能够将动态图数据转化为鲁棒的动态节点嵌入。

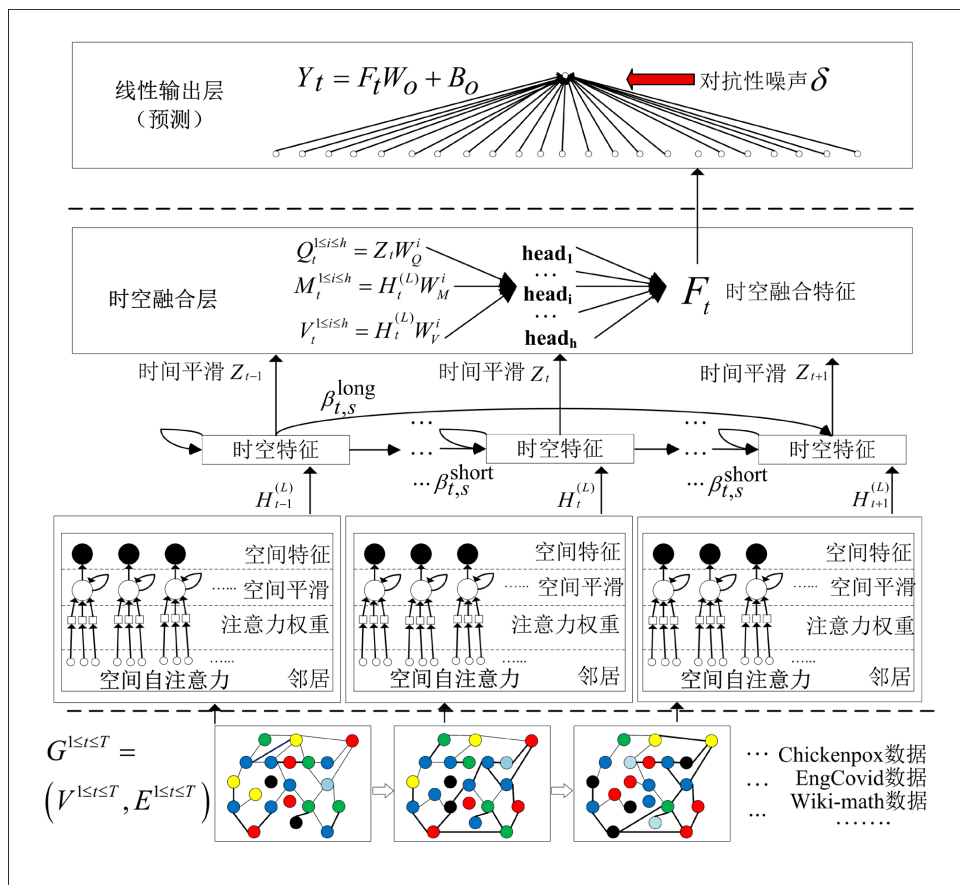


Figure 1. Robust graph spatial-temporal attention network against projected gradient attacks  
图 1. 防御投影梯度攻击的鲁棒图时空注意力网络

## 4.1. 投影梯度下降对抗攻击

### 4.1.1. 对抗攻击与防御

面向图时空模型的对抗攻击通过扰动训练集中各快照的节点与边来干扰模型训练。面向结构的对抗攻击可以修改图的拓补结构，使得两个完全无关的节点间存在连接，导致邻接矩阵发生了彻底的变化。对节点的特征或者标签进行对抗攻击，能使模型对节点产生错误的判断。

攻击者恶意更改图神经网络模型的训练集  $G^{1 \sim T}$  每个快照的节点或边，以降低模型的精度，防御者需要抵抗训练集  $G^{1 \sim T}$  被更改所带来的性能下降。应对对抗攻击的模型鲁棒性指当数据被恶意更改时，模型依旧能保持其性能。故而防御者的目标是找到最优模型  $\hat{f}$ ，最小化攻击者的干扰，使得模型在测试集上依旧能够最小化模型输出与标签间的差距：

$$\arg \min_{\hat{f}} \mathbb{E}_{1 \leq t \leq T} \text{Loss}_{\text{RSTGNN}} \left( \hat{f}_{X', A'} \left( X_{\text{test}}^t, A_{\text{test}}^t \right), Y_{\text{test}}^t \right) \quad (5)$$

$$\text{s.t.} : \arg \min_{X', A'} \text{Loss}_{\text{atk}} \left( f_{X', A'} \left( X_{\text{test}}, A_{\text{test}} \right), Y_{\text{test}} \right) \quad (6)$$

### 4.1.2. 攻击目标

投影梯度下降法(Projected Gradient Descent, PGD)是一种常用的迭代式对抗攻击方法,通过在允许的扰动范围内沿着损失函数梯度方向不断更新输入,生成能最大化模型损失的对抗样本。在图时空网络中,PGD攻击者在原始图结构或节点特征的基础上添加一个小扰动 $\delta$ ,并在每一次迭代中根据损失函数关于输入的梯度方向更新扰动。通过多次迭代更新,PGD攻击能够在给定的扰动约束范围(如 $\|\delta\|_p \leq \varepsilon$ )内找到使模型损失最大的扰动,从而得到最具攻击性的对抗样本。

$$\delta^{k+1} = \Pi_S \left( \delta^k + \zeta \cdot \text{sign} \left( \nabla_{X,A} L_{atk} \left( f \left( X + \delta_X^k, A + \delta_A^k \right), Y \right) \right) \right) \quad (7)$$

其中, $\zeta$ 表示每一步的更新步长, $\nabla_{X,A}$ 表示损失函数对节点特征或邻接矩阵的梯度, $\Pi_S(\cdot)$ 表示投影操作,用于将更新后的扰动限制在允许的扰动集合 $S$ 内, $\delta_X$ 与 $\delta_A$ 分别表示对节点特征和图结构的扰动。

### 4.2. 强化空间层

在图时空网络模型中,节点特征的提取依赖于邻域信息的有效聚合,PGD对抗攻击下会干扰空间信息传播。为提升模型在空间特征提取上的鲁棒性,本文在图注意力网络的基础上引入邻域平滑约束。故而快照 $t$ 在第 $l+1$ 层的表示更新为:

$$H_t^{(l+1)} = \sigma \left( \frac{1}{K} \sum_{k=1}^K \alpha_t^{(k)} H_t^{(l)} W^{(k)} \right) \quad (8)$$

$$\alpha_{t,ij}^{(k)} = \frac{\exp \left( \text{LeakyReLU} \left( \mathbf{a}^{(k)} \cdot \left[ W^{(k)} h_{t,i} \parallel W^{(k)} h_{t,j} \right] \right) \right)}{\sum_{j \in N_i} \exp \left( \text{LeakyReLU} \left( \mathbf{a}^{(k)} \cdot \left[ W^{(k)} h_{t,i} \parallel W^{(k)} h_{t,j} \right] \right) \right)} \quad (9)$$

其中, $N_i$ 表示节点 $i$ 的邻居集合, $K$ 表示注意力头数量, $W^{(k)} \in \mathbb{R}^{d \times d}$ 为第 $k$ 个注意力头的可学习权重矩阵, $\sigma(\cdot)$ 表示非线性激活函数, $\parallel$ 表示向量拼接操作。节点特征归一化的注意力矩阵 $\alpha_t^{(k)}$ 由元素 $\alpha_{t,ij}^{(k)} \in \mathbb{R}$ 构成,包含了节点 $i$ 从邻居 $j$ 接收信息的权重, $\mathbf{a}^{(k)} \in \mathbb{R}^{n \times n}$ 是第 $k$ 个注意力头的注意力权重向量。

为了降低对抗扰动带来的异常信息传播,该层引入空间平滑机制,该机制融合节点特征与其领域的平均特征用以增强节点局部的一致性约束:

$$\tilde{H}_t^{(L)} = (1 - \lambda) H_t^{(L)} + \lambda D_t^{-1} A_t H_t^{(L)} \quad (10)$$

其中, $D_t^{-1}$ 为快照 $t$ 的度矩阵, $\lambda$ 为空间平滑系数,该系数用于调控节点特征与领域信息之间的权衡,该平滑机制能够缓解由于对抗扰动导致的局部特征偏移问题,从而增强模型在空间结构上的稳定性。 $\lambda$ 越小则平滑效果越低,难以抵抗PGD扰动,如果 $\lambda$ 过大则会过于削弱节点自身特征表示,故而 $\lambda$ 的取值范围需要偏向于节点自身特征,可取值 $\lambda \in [0.1, 0.3]$ 之间。

### 4.3. 强化时间层

图时空网络从动态图数据中提取的节点特征会随时间演化,而对抗PGD扰动会破坏数据中时间依赖关系,从而影响模型对时空模式的学习能力。因此,时间网络层将时间注意力机制与时间一致性约束相结合,以提升模型对时间模式扰动的鲁棒性。强化时间层结合了长短双尺度自注意力机制,能够在时间维度上实现节点特征聚合:

$$Z_t^{\text{long}} = \sum_{s=1}^{T'} \beta_{t,s}^{\text{long}} \tilde{H}_s^{(L)}; Z_t^{\text{short}} = \sum_{s=t-\tau+1}^t \beta_{t,s}^{\text{short}} \tilde{H}_s^{(L)} \quad (11)$$

$$\beta_{t,s}^{\text{long}} = \frac{\exp\left(f\left(\tilde{H}_t^{(L)}, \tilde{H}_s^{(L)}\right)\right)}{\sum_{r=1}^{T'} \exp\left(f\left(\tilde{H}_t^{(L)}, \tilde{H}_r^{(L)}\right)\right)}; \beta_{t,s}^{\text{short}} = \frac{\exp\left(f\left(\tilde{H}_t^{(L)}, \tilde{H}_s^{(L)}\right)\right)}{\sum_{r=t-\tau+1}^t \exp\left(f\left(\tilde{H}_t^{(L)}, \tilde{H}_r^{(L)}\right)\right)} \quad (12)$$

其中,  $T'$  和  $\tau$  分别表示长时间窗口和短时间窗口长度,  $H_s^{(L)}$  表示由空间层所提取的快照  $s$  的节点表示, 长短时间注意力权重  $\beta_{t,s}^{\text{long}}$  与  $\beta_{t,s}^{\text{short}}$  用于表示快照  $t$  对于快照  $s$  的关注程度。

当历史快照受到扰动时, 其与当前时刻的相关性下降, 其注意力权重在分配中被抑制。故而, 长时间注意力用于提供稳定的全局上下文, 而短时间注意力能够增强对局部动态变化的敏感性。为避免固定加权方式带来的信息偏置, 强化时间层引入门控机制对长短期信息进行自适应融合:

$$U_t = [Z_t^{\text{long}} \parallel Z_t^{\text{short}}]; \mu_t = \sigma(U_t W_g + b_g) \quad (13)$$

$$Z_t = \mu_t \odot Z_t^{\text{long}} + (1 - \mu_t) \odot Z_t^{\text{short}} \quad (14)$$

其中,  $W_g \in \mathbb{R}^{2d \times d}$  为降维用的权重矩阵,  $b_g \in \mathbb{R}^d$  为偏置项,  $\mu_t \in \mathbb{R}^{n \times d}$  能够自适应调控长短期信息的贡献, 从而缓解对抗扰动对特定快照的影响。

当动态图中的部分快照受到对抗扰动时, 不同历史快照对当前时刻特征提取的重要性也不一样。如果只是平均地分配时间注意力权重, 则意味着默认所有前置快照都是同等可靠的, 攻击者就容易在部分历史快照中引入对抗扰动。故而, 时间层建立了相似度函数  $f(\cdot)$  以衡量快照  $t$  与快照  $s$  的相关性, 并据此适应性地分配时间自注意力权重。如果某快照  $s$  的节点特征被扰动, 则与当前快照  $t$  的一致性会降低, 其被分配的时间注意力权重也更少, 从而缓解异常历史快照的扰动:

$$f\left(\tilde{H}_t^{(L)}, \tilde{H}_s^{(L)}\right) = \mathbf{v}^\top \tanh\left(W_q \tilde{H}_t^{(L)} + W_k \tilde{H}_s^{(L)}\right) \quad (15)$$

其中,  $W_q \in \mathbb{R}^{d \times d}$  和  $W_k \in \mathbb{R}^{d \times d}$  为可学习权重矩阵,  $\mathbf{v}$  为时间注意力向量。

另外, 时间网络层也通过时间平滑约束来减少对抗扰动引起的时间特征异常:

$$\tilde{Z}_t = \gamma Z_t + (1 - \gamma) Z_{t-1} \quad (16)$$

其中, 时间平滑系数  $\gamma \in [0.1, 0.5]$  调控当前快照与历史快照在时间特征上的权衡, 若  $\gamma$  过大, 当前的节点表示会过于依赖历史快照。对比空间平滑系数, 节点间关系由相对稳定的图结构所决定, 领域信息的一致性较强, 通过较小程度的平滑就能抑制局部扰动。而在时间维度上, 引入更大的时间平滑系数能够适当增强平滑强度来抵消瞬时扰动带来的异常。

#### 4.4. 时空融合层

在 RGSTAN 的预测任务中, 时空两个网络层是将空间与时间特征分开处理, 缺乏对时空特征关联的建模, 故而构建了时空融合层, 使用交叉自注意力将动态节点的原始特征、平滑空间特征与平滑时间特征进行动态加权融合, 使模型可以在不同快照上适应性地关注空间模式或时间模式, 能够避免由局部过度扰动带来空间模式受损, 也能降低少数节点被长期影响所带来的错误演化, 从而提升 RGSTAN 模型的鲁棒性。时空融合层的特征表示为:

$$Q_t = \tilde{Z}_t W_Q, \quad M_t = H_t^{(L)} W_M, \quad V_t = \tilde{H}_t^{(L)} W_V \quad (17)$$

$$F_t = \tilde{H}_t^{(L)} + \text{MultiHead}(Q_t, M_t, V_t) \quad (18)$$

其中,  $W_Q \in \mathbb{R}^{d \times d}$ ,  $W_M \in \mathbb{R}^{d \times d}$ ,  $W_V \in \mathbb{R}^{d \times d}$  为时空融合层的权重矩阵,  $Q_t \in \mathbb{R}^{n \times d}$ 、 $M_t \in \mathbb{R}^{n \times d}$ 、 $V_t \in \mathbb{R}^{n \times d}$  分别表示时空融合层中交叉注意力的查询、键和值矩阵,  $M_t$  由未通过空间平滑的原始特征映射而来, 可以将节点原始特征融入时空模式中。

$$\text{head}_i = \text{softmax} \left( \frac{Q_t^i (M_t^i)^\top}{\sqrt{d/h}} \right) V_t^i \quad (19)$$

$$\text{MultiHead}(Q_t, M_t, V_t) = \text{Concat}(\text{head}_1, \dots, \text{head}_h) W^R \quad (20)$$

其中,  $\text{head}_i$  表示该层的第  $i$  个注意力, 基于交叉注意力的时空融合层将多个注意力聚合形成更加丰富的时空特征表示。该过程不但实现不同时空特征的有效聚合, 也提升模型对复杂时空动态的建模能力与鲁棒性。

#### 4.5. 节点回归损失函数

本文提出的 RGSTAN 模型主要用于节点回归任务, 即对未来的节点状况进行预测, 该模型的优化目标  $L_{\text{reg}}$  可通过 MSE 指标衡量。当时空融合层的输出为  $F_t$  时, 模型最终的线性层能够将节点的  $F_t$  特征从  $n \times d$  映射到  $n \times 1$  以完成预测, 其输出形式为:

$$\hat{Y}_t = F_t W_o + b_o \quad (21)$$

其中,  $W_o \in \mathbb{R}^{d \times 1}$  和  $b_o \in \mathbb{R}^1$  分别为可学习的权重矩阵和偏置项。

故而, 执行节点回归任务的 RGSTAN 的损失函数为:

$$L_{\text{reg}} = \frac{1}{TN} \sum_{t=1}^T \sum_{i=1}^N (y_i^t - \hat{y}_i^t)^2 \quad (22)$$

为了优化模型对于 PGD 攻击的鲁棒性, RGSTAN 在损失函数上添加了对抗扰动, 使得最后的模型优化目标 Loss 为:

$$L_{\text{rob}} = \frac{1}{TN} \sum_{t=1}^T \sum_{i=1}^n (y_i^t - \hat{y}_i^t(X + \delta))^2 \quad (23)$$

$$\text{Loss} = L_{\text{reg}} + \eta L_{\text{rob}} \quad (24)$$

其中,  $\eta$  为权衡参数, 用于控制节点回归损失与对抗扰动损失之间的比例。通过在训练过程中引入对抗样本, 能够降低模型对于输入扰动的敏感性, 从而在扰动环境下学习更加鲁棒的时空特征表示。

## 5. 实验设计

### 5.1. 实验设定

1) 实验环境与参数设定: 本实验采用 NVIDIA RTX 4070 GPU 执行, 其中 CUDA 版本为 11.8。另外, 采用 PyTorch 2.3.1 + Python 3.8 模拟 PGD 攻击并实现模型构建与训练。图时空模型通过 Adam 优化器进行模型参数更新, 其中学习率为  $1 \times 10^{-3}$ , 衰减系数设为  $5 \times 10^{-4}$ , 训练轮数 epoch 设置为 200。实验用动态图数据集按照时间跨度进行划分, 前 60% 的图快照作为训练集, 后 40% 的图快照作为测试集。此外, 模型采用的注意力头数  $K$  设为 8, 空间平滑系数  $\lambda$  取 0.15, 时间平滑系数  $\gamma$  取 0.25, 权衡参数  $\eta$  为 0.2, 时间窗口长度  $T'$  与  $\tau$  设为 4。在对抗训练中, PGD 扰动强度  $\varepsilon$  设为 (0.0, 0.05, 0.2), 迭代步数为  $\Lambda = 5$ , 步长按照  $\zeta = \varepsilon/\Lambda$  设置。为全面评估模型的预测性能与鲁棒性, 本文选取平均绝对误差(MAE)和均方根误差(RMSE)作为回归任务的评价指标以衡量模型预测精度。

2) 动态图数据: 本实验所采取的动态图数据[4]为 Chickenpox、England-Covid、Wiki-Math(见表 1)。其中, Chickenpox 以匈牙利各地区之间的关联为节点的邻接关系, 节点特征表示为各地区每周的水痘病例数。EngCovid 数据集以英国各地区为节点, 通过各地区间的疫情人口流动关系构建节点特征与节点关联。Wiki-Math 数据集则以维基百科页面为节点, 其具有超链接关系为边, 其中节点特征为随时间变化的用户访问量。图时空模型能够学习动态图数据在空间上的传播模式和时间上的动态规律, 从而预测下一个图快照的节点特征, 进而推理出疫情变化、网络流量。

**Table 1.** Dynamic graph data

**表 1.** 动态图数据

名称	节点数	边数	时间步	采样	特征维度
Chickenpox	20	102	517	每日	4
EngCovid	129	2158	53	每日	8
Wiki-Math	1068	27,079	723	每日	8

3) 图时空模型: 本实验通过三种防御方法(AT、DD、IP)在四类图时空网络(见表 2)中抵御不同强度的 PGD 攻击以证明所提出 RGSTAN 模型的鲁棒性。

**Table 2.** Graph spatial-temporal models

**表 2.** 图时空模型

模型	模型构成	特点
AGCRN [14]	自适应图卷积 + 门控循环	自适应学习隐藏图结构
GCLSTM [15]	图卷积嵌入 LSTM 单元	联合建模空间与时间依赖
GConvGRU [16]	以图卷积替代 GRU 全连接结构	能显式建模图结构中的空间依赖
TGCN [17]	图卷积 + GRU 紧凑组合	结构简洁、训练稳定

## 5.2. 实验分析

本实验模拟了在不同程度投影梯度攻击下图时空网络模型的预测精度, 本文所提出的 RSTGN 模型采用时空平滑 STS (Spatial-Temporal Smoothing)方法以抵御 PGD 攻击, 其余模型分别采用对抗攻击 AT, 防御蒸馏 DD 以及输入净化 IP 以提升其鲁棒性。

如表 3 所示, 在无扰动  $\varepsilon = 0$  下, 各模型整体表现较为接近, 但不同防御策略之间仍存在差异。对于 AGCRN、GCLSTM 和 GConvGRU 模型, AT、DD 与 IP 三种策略在不同数据集上的效果波动较大, 例如在 EngCovid 数据集上, IP 策略虽然在部分情况下可以降低 MAE (如 AGCRN 降至 0.781), 但 RMSE 为 1.022 仍存在不稳定现象。DD 策略在 Wiki-Math 数据集上反而导致误差上升(如 GCLSTM 的 RMSE 达到 0.743), 其高于 AT (0.714)。AT 采取的扰动样本训练会干扰原始数据分布, DD 会削弱时空特征的建模能力, IP 容易破坏原有的时空结构从而导致结果不稳定。对比来看, RGSTAN 在空间与时间两个维度上同时抑制噪声传播, 采用 STS 策略的 RGSTAN 在三个数据集上均保持较低误差(如 RMSE 分别为 0.986、0.939 和 0.704), 说明其在无攻击情况下已具备良好的泛化能力。表 3 提供了面向 RGSTAN 的消融实验结果, 其中时空平滑权衡参数被分别设置为最大取值, 以及单方面的平滑从而证明时空平滑策略对模型性能的影响。当时空平滑程度达到最大, 过强的平滑导致模型在三个数据集上的误差有所上升。另外, 当模型除去时间平滑后, 模型性能的下降比除去空间平滑更严重, 表明时间维度信息对于动态图数据的

建模影响更大。该消融时间表明 RGSTAN 能够通过合理调控时空平滑参数来稳固其时空依赖，从而维持模型的精度。

**Table 3.** Accuracy of different graph spatial-temporal models under perturbation  $\varepsilon = 0$

**表 3.** 扰动强度  $\varepsilon = 0$  下的各类图时空模型精度

模型	防御	指标	Chickenpox	EngCovid	Wiki-Math
AGCRN	AT	MAE	0.683	1.067	0.507
		RMSE	1.029	1.393	0.801
	DD	MAE	0.666	0.960	0.563
		RMSE	1.013	1.238	0.849
	IP	MAE	0.660	0.781	0.559
		RMSE	1.002	1.022	0.858
GCLSTM	AT	MAE	0.676	0.924	0.419
		RMSE	0.991	1.226	0.714
	DD	MAE	0.647	0.764	0.446
		RMSE	0.975	1.030	0.743
	IP	MAE	0.672	1.009	0.538
		RMSE	0.992	1.313	0.829
GConvGRU	AT	MAE	0.771	0.948	0.426
		RMSE	1.090	1.249	0.721
	DD	MAE	0.662	0.706	0.462
		RMSE	0.978	0.952	0.746
	IP	MAE	0.672	0.977	0.534
		RMSE	0.987	1.281	0.824
TGCN	AT	MAE	0.700	0.894	0.413
		RMSE	1.036	1.182	0.709
	DD	MAE	0.664	0.861	0.615
		RMSE	0.994	1.153	0.929
	IP	MAE	0.710	0.972	0.625
		RMSE	1.035	1.270	0.937
RGSTAN (STS)	$\lambda = 0.15; \gamma = 0.25$	MAE	0.656	0.702	0.409
		RMSE	0.986	0.939	0.704
	$\lambda = 0.3; \gamma = 0.5$	MAE	0.684	0.771	0.459
		RMSE	1.028	1.036	0.742
	$\lambda = 0.15$ ; 无时间平滑	MAE	0.706	0.842	0.479
		RMSE	1.051	1.087	0.791
	$\gamma = 0.25$ ; 无空间平滑	MAE	0.693	0.807	0.466
		RMSE	1.037	1.063	0.756

如表 4 所示,  $\varepsilon = 0.05$  时各模型性能开始明显下降, 且不同防御方法的鲁棒性差异进一步放大。传统模型在 IP 策略下普遍出现较大波动, 例如 GCLSTM 在 EngCovid 数据集上的 RMSE 升至 1.322, 显示出对 PGD 攻击的敏感性; 而 DD 策略虽在部分数据(如 GConv-GRU 在 EngCovid 上的 RMSE 为 0.970)取得一定改善, 但整体表现不稳定。RGSTAN 在该扰动强度下依然保持显著优势, 在图时空模型中, PGD 扰动会通过邻接关系进行传播, 并通过图演化在时间上不断累积, 而时空平滑机制能够在传播过程中持续削弱噪声信号, 从而抑制误差放大。在 Wiki-Math 数据集上的 MAE 和 RMSE 分别达到 0.405 和 0.699, 远低于其他模型, 同时在 EngCovid 数据集上 RMSE 仅为 0.934, 表明 RGSTAN 在中等强度攻击下仍能有效保持模型精度。根据表 4 中的消融实验结果可知, 不同的时空平滑策略在 PGD 扰动下对模型鲁棒性的影响更明显; 过强的时空平滑约束虽然能抑制对抗攻击, 但也会对模型本身的性能造成干扰。在表 4 中, RGSTAN 不采用时间平滑的影响(不同数据集上 RMSE 为 1.053、1.126、0.798)依旧高于不采用空间平滑(1.036、1.061、0.768), 这是因为空间结构的扰动主要是影响的当前局部图结构, 而时间维度上的干扰会随着图演化而累积。

**Table 4.** Accuracy of different graph spatial-temporal models under perturbation  $\varepsilon = 0.05$

**表 4.** 扰动强度  $\varepsilon = 0.05$  下的各类图时空模型精度

模型	防御	指标	Chickenpox	EngCovid	Wiki-Math
AGCRN	AT	MAE	0.704	0.952	0.548
		RMSE	1.018	1.261	0.837
	DD	MAE	0.666	0.978	0.616
		RMSE	1.014	1.230	0.903
	IP	MAE	0.659	0.806	0.573
		RMSE	1.011	1.049	0.866
GCLSTM	AT	MAE	0.718	0.981	0.538
		RMSE	1.037	1.286	0.809
	DD	MAE	0.668	0.793	0.627
		RMSE	1.000	1.040	0.886
	IP	MAE	0.679	1.019	0.563
		RMSE	1.001	1.322	0.854
GConvGRU	AT	MAE	0.736	0.967	0.546
		RMSE	1.030	1.268	0.817
	DD	MAE	0.704	0.750	0.558
		RMSE	1.016	0.970	0.826
	IP	MAE	0.686	0.993	0.566
		RMSE	1.024	1.295	0.855
TGCN	AT	MAE	0.702	0.857	0.658
		RMSE	1.019	1.140	0.968
	DD	MAE	0.654	0.896	0.695
		RMSE	1.010	1.183	0.986
	IP	MAE	0.714	0.985	0.675
		RMSE	1.055	1.281	0.998

续表

<b>RGSTAN (STS)</b>	$\lambda = 0.15; \gamma = 0.25$	MAE	0.651	0.699	0.405
		RMSE	0.982	0.934	0.699
	$\lambda = 0.3; \gamma = 0.5$	MAE	0.679	0.754	0.444
		RMSE	1.018	0.997	0.739
	$\lambda = 0.15$ ; 无时间平滑	MAE	0.709	0.873	0.488
		RMSE	1.053	1.126	0.798
	$\gamma = 0.25$ ; 无空间平滑	MAE	0.695	0.819	0.471
		RMSE	1.036	1.061	0.768

如表 5 所示, 扰动强度  $\varepsilon = 0.2$  时, 模型间的性能差距进一步扩大。AGCRN、GCLSTM 及 TGCN 等模型在不同防御策略下均出现不同程度的性能退化, 尤其是在 Wiki-Math 数据集上, RMSE 普遍接近或超过 0.9, 部分甚至超过 1.0 (如 TGCN 在 IP 策略下达到 1.008)。同时, 各防御方法之间仍缺乏一致性, IP 和 DD 策略在不同数据集上的效果差异明显。相比之下, RGSTAN 在三个数据集上依然保持稳定表现, 其 MAE 与 RMSE 均处于较低水平(如 Wiki-Math 上 MAE 为 0.564, RMSE 为 0.883), 体现出较强的抗扰动能力。表 5 中的消融实验进一步说明了随着 PGD 扰动程度增强, 扰动在时间维度上的累积也更大, 从而干扰了模型性能。然而, 不同平滑配置上的性能变化也受 PGD 本身攻击策略的影响, 因为 PGD 在扰动动态图的各个图快照时具有不均匀性, 部分区域会有随机干扰的效果。例如, 在 EngCovid 数据集上, 时空平滑参数的变化对模型的鲁棒性影响较大。

**Table 5.** Accuracy of different graph spatial-temporal models under perturbation  $\varepsilon = 0.2$

**表 5.** 扰动强度  $\varepsilon = 0.2$  下的各类模型精度

模型	防御	指标	Chickenpox	EngCovid	Wiki-Math
<b>AGCRN</b>	<b>AT</b>	MAE	0.723	0.980	0.578
		RMSE	1.032	1.292	0.865
	<b>DD</b>	MAE	0.680	0.986	0.634
		RMSE	1.025	1.239	0.919
	<b>IP</b>	MAE	0.674	0.832	0.588
		RMSE	1.023	1.076	0.880
<b>GCLSTM</b>	<b>AT</b>	MAE	0.726	0.995	0.576
		RMSE	1.042	1.297	0.817
	<b>DD</b>	MAE	0.681	0.804	0.658
		RMSE	1.008	1.048	0.905
	<b>IP</b>	MAE	0.695	1.029	0.571
		RMSE	1.011	1.331	0.861
<b>GConvGRU</b>	<b>AT</b>	MAE	0.755	0.980	0.580
		RMSE	1.042	1.278	0.841
	<b>DD</b>	MAE	0.721	0.762	0.587
		RMSE	1.027	0.979	0.845
	<b>IP</b>	MAE	0.712	1.007	0.574
		RMSE	1.042	1.307	0.862

续表

TGCN	AT	MAE	0.722	0.877	0.685
		RMSE	1.034	1.156	0.994
	DD	MAE	0.660	0.905	0.708
		RMSE	1.016	1.191	0.997
	IP	MAE	0.761	0.997	0.685
		RMSE	1.085	1.292	1.008
RGSTAN (STS)	$\lambda = 0.15; \gamma = 0.25$	MAE	0.669	0.863	0.564
		RMSE	1.003	1.118	0.883
	$\lambda = 0.3; \gamma = 0.5$	MAE	0.709	0.918	0.586
		RMSE	1.052	1.173	0.907
	$\lambda = 0.15$ ; 无时间平滑	MAE	0.734	0.947	0.618
		RMSE	1.086	1.221	0.963
$\gamma = 0.25$ ; 无空间平滑	MAE	0.718	0.934	0.602	
	RMSE	1.068	1.205	0.929	

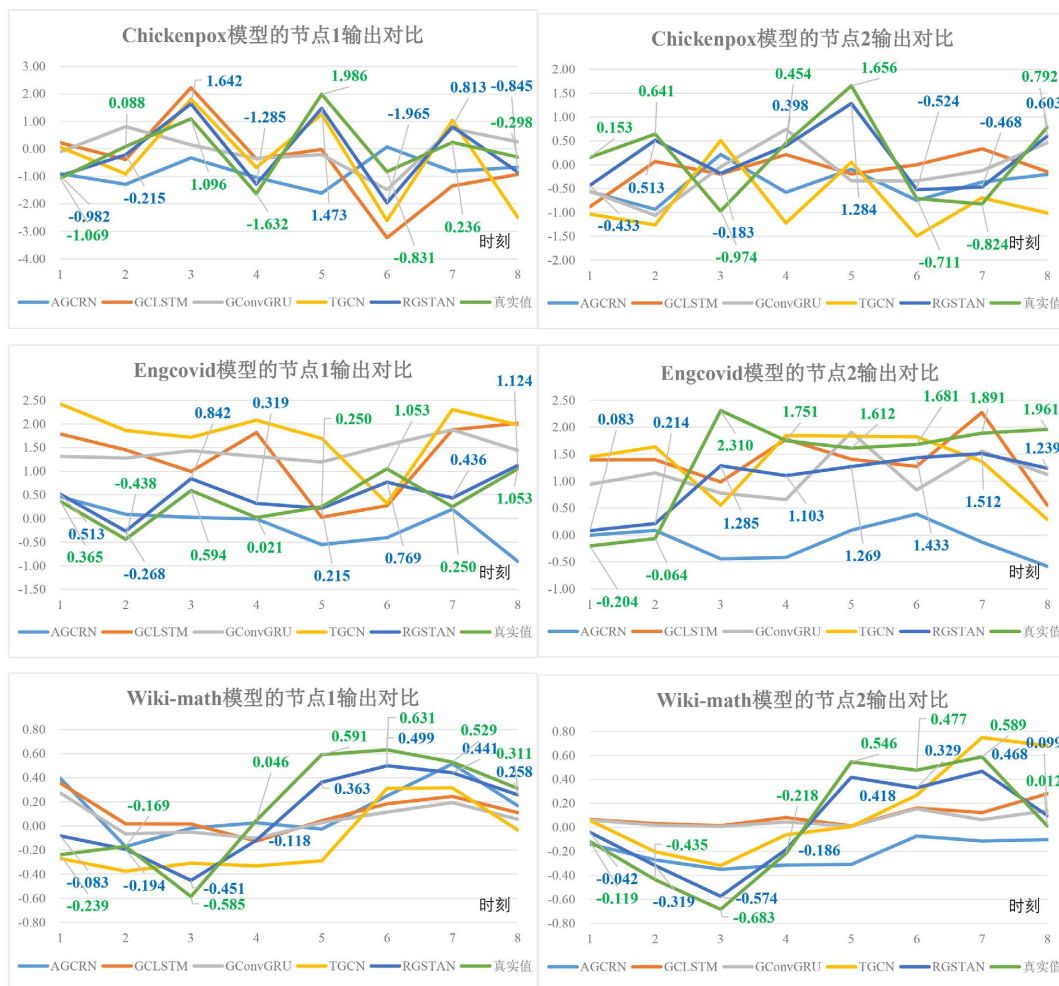


Figure 2. Comparison of model outputs on two high-degree nodes under high PGD perturbations

图 2. 高 PGD 扰动下不同模型在两个高度数节点上的输出对比

由三种扰动约束下的实验结果可知, AT、DD 或 IP 等策略随着 PGD 攻击强度的增加, 难以在不同数据集上保持稳定性能。而采用时空平滑的 RGSTAN 模型在三种条件下能够维持更小的性能波动, 尤其在中高强度攻击下优势更加显著。这表明时空平滑 STS 机制能够有效缓解对抗扰动对图结构与图演化信息的破坏, 从而实现更鲁棒的时空预测。另外, 基于对不同时空平滑参数的消融实验表明, RGSTAN 对于时空平滑参数的敏感性较为稳定, 在  $\lambda = 0.15; \gamma = 0.25$  附近时, 模型能够有效保持其对 PGD 攻击的鲁棒性, 但当平滑程度达到最高时也出现了过平滑问题, 从而造成模型性能下滑。而且从表 3~5 可知, 图时空模型 RGSTAN 对于时间平滑参数具有更强的敏感性。

在高强度 PGD 扰动下, 本实验基于三个数据中最高度数的两个节点在 8 个时刻上的预测结果(见图 2), 通过计算各模型输出与真实值之间的由时刻 1 到时刻 8 的平均误差(MAE)来表示模型是否能实现高度节点的鲁棒时空建模。在节点上的结果表明, RSTGAN ( $\lambda = 0.15; \gamma = 0.25$ )在 Chickenpox 数据上的平均误差为 0.41 和 0.60, 明显低于 GCLSTM (1.10 和 0.81)与 TGCN (1.31 和 0.86), 接近于 AGCRN (0.45 与 0.57)。在 EngCovid 上, RGSTAN 的平均误差为 0.3 和 0.56, 其误差优于其余四个模型。在 Wiki-Math 上, RGSTAN 拥有比其他模型更低的平均误差。在 Chickenpox 与 EngCovid 的部分时刻, RGSTAN 也出现过明显的预测偏移, 这类异常主要来自于高强度 PGD 对局部时空依赖的破坏。由于高连接度的节点聚合了更多领域信息, 其在受到扰动后更容易引入噪声。但对比其他模型, RGSTAN 的误差波动幅度整体更小, 且能在后续的时刻降低与真实值的差距, 说明 RSTGAN 在高强度 PGD 扰动下具有更好的鲁棒性。

## 6. 结论

本文提出了基于时空平滑的鲁棒图时空注意力网络模型, 该模型采用交叉时空自注意力实现了图结构信息与图演化信息的联合建模, 克服时空特征分离的表征局限也增强了其动态模式的捕获能力。另外, 该模型通过时空平滑策略, 从时间维度和空间维度实施跨时刻或邻域聚合扰动的抑制。本文在三个真实动态图数据集上展开了 PGD 攻击实验以表明 RGSTAN 的鲁棒性。本研究也为构建智能交通、疾病预警等关键鲁棒系统构建提供了理论支撑, 未来将进一步探究图时空网络模型的鲁棒性与可解释性, 从而支撑图时空模型在安全敏感场景中实际部署。

## 基金项目

重庆电子科技职业大学 2025 年学生科技创新“马跃”工程项目资助(课题批准号: 25XJXSCX03)。

## 参考文献

- [1] Leskovec, J. and Krevl, A. (2014) SNAP Datasets: Stanford Large Network Dataset Collection. <http://snap.stanford.edu/data>
- [2] 乔少杰, 薛骐, 杨国平, 等. 基于动态自适应时空图的多元时序预测模型[J]. 计算机学报, 2024, 47(12): 2925-2937.
- [3] 何玉林, 赖俊龙, 崔来中, 等. 基于时空注意力的多粒度链路预测算法[J]. 软件学报, 2025, 36(9): 4311-4326.
- [4] Rozemberczki, B., Scherer, P., He, Y., Panagopoulos, G., Riedel, A., Astefanoaei, M., et al. (2021) Pytorch Geometric Temporal: Spatiotemporal Signal Processing with Neural Machine Learning Models. *Proceedings of the 30th ACM International Conference on Information & Knowledge Management*, Virtual Event, 1-5 November 2021, 4564-4573. <https://doi.org/10.1145/3459637.3482014>
- [5] Madry, A., Makelov, A., Schmidt, L., Tsipras, D. and Vladu, A. (2018) Towards Deep Learning Models Resistant to Adversarial Attacks. arXiv:1706.06083.
- [6] 先兴平, 吴涛, 乔少杰, 等. 图学习隐私与安全问题研究综述[J]. 计算机学报, 2023, 46(6): 1184-1212.
- [7] 金柯君, 于洪涛, 吴翼腾, 等. 基于改进投影梯度下降算法的图卷积网络投毒攻击[J]. 计算机工程, 2022, 48(10): 176-183.

- 
- [8] 柏杨, 陈晋音, 郑海斌, 等. 面向图垂直联邦学习的对抗攻击方法[J]. 计算机科学, 2025, 52(S2): 841-850.
- [9] Li, Y., Jin, W., Xu, H. and Tang, J. (2020) Deeprobust: A Pytorch Library for Adversarial Attacks and Defenses. arXiv:2005.06149.
- [10] Zi, B., Zhao, S., Ma, X. and Jiang, Y. (2021) Revisiting Adversarial Robustness Distillation: Robust Soft Labels Make Student Better. 2021 *IEEE/CVF International Conference on Computer Vision (ICCV)*, Montreal, 10-17 October 2021, 16443-16452. <https://doi.org/10.1109/iccv48922.2021.01613>
- [11] Lee, W. and Park, H. (2025) Self-Supervised Adversarial Purification for Graph Neural Networks. *Proceedings of the 42nd International Conference on Machine Learning*, Vancouver, 13-19 July 2025, 33715-33735.
- [12] 王煜恒, 刘强, 伍晓洁. RCGNN: 图注入攻击下的图神经网络鲁棒性认证方法[J]. 计算机工程与科学, 2025, 47(3): 434-447.
- [13] 王新哲, 孙望舒, 罗晨, 等. 基于动态时空图网络的数据安全态势预警技术[J/OL]. 计算机与现代化, 1-15. <https://link.cnki.net/urlid/36.1137.tp.20251110.1806.007>, 2026-03-15.
- [14] Bai, L., Yao, L., Li, C., et al. (2020) Adaptive Graph Convolutional Recurrent Network for Traffic Forecasting. *Advances in Neural Information Processing Systems*, **33**, 17804-17815.
- [15] Chen, J., Wang, X. and Xu, X. (2022) GC-LSTM: Graph Convolution Embedded LSTM for Dynamic Network Link Prediction. *Applied Intelligence*, **52**, 7513-7528. <https://doi.org/10.1007/s10489-021-02518-9>
- [16] Seo, Y., Defferrard, M., Vandergheynst, P. and Bresson, X. (2018) Structured Sequence Modeling with Graph Convolutional Recurrent Networks. In: Cheng, L., Leung, A. and Ozawa, S., Eds., *Lecture Notes in Computer Science*, Springer International Publishing, 362-373. [https://doi.org/10.1007/978-3-030-04167-0\\_33](https://doi.org/10.1007/978-3-030-04167-0_33)
- [17] Zhao, L., Song, Y., Zhang, C., Liu, Y., Wang, P., Lin, T., et al. (2020) T-GCN: A Temporal Graph Convolutional Network for Traffic Prediction. *IEEE Transactions on Intelligent Transportation Systems*, **21**, 3848-3858. <https://doi.org/10.1109/tits.2019.2935152>