
Derivation of Back-propagation for Graph Convolutional Networks using Matrix Calculus and its Application to Explainable Artificial Intelligence

Yen-Che Hsiao

Department of Electrical and Computer Engineering
University of Connecticut
Storrs CT 06269, USA
yen-che.hsiao@uconn.edu

Rongting Yue

Department of Electrical and Computer Engineering
University of Connecticut
Storrs CT 06269, USA

Abhishek Dutta

Department of Electrical and Computer Engineering
University of Connecticut
Storrs CT 06269, USA

Abstract

This paper provides a comprehensive and detailed derivation of the backpropagation algorithm for graph convolutional neural networks using matrix calculus. The derivation is extended to include arbitrary element-wise activation functions and an arbitrary number of layers. The study addresses two fundamental problems, namely node classification and link prediction. To validate our method, we compare it with reverse-mode automatic differentiation. The experimental results demonstrate that the median sum of squared errors of the updated weight matrices, when comparing our method to the approach using reverse-mode automatic differentiation, falls within the range of 10^{-18} to 10^{-14} . These outcomes are obtained from conducting experiments on a five-layer graph convolutional network, applied to a node classification problem on Zachary's karate club social network and a link prediction problem on a drug-drug interaction network. Finally, we show how the derived closed-form solution can facilitate the development of explainable AI and sensitivity analysis.

1 Introduction

Graph neural network (GNN) is a neural network model for processing data represented in graph domains, encompassing cyclic, directed, and undirected graphs [1]. Graph convolutional network (GCN) is a type of GNN model that employs layer-wise propagation rule, operating directly on graphs [2]. A GCN layer comprises message passing over nodes/edges, succeeded by an aggregation/pooling strategy and a fully connected layer [3].

Node classification is one of the most common research directions in graph analysis, where the objective of the task is to predict a specific class for each unlabeled node in a graph using graph information

[4]. J. Zhang *et al.* [5] applied GCN for distribution system anomaly detection, categorizing nodes into normal data, data anomalies, or event anomalies, and implemented their method on the IEEE 37-node distribution systems. Link prediction seeks to infer unobserved/missing links or predict future ones based on the connections within currently observed partial networks [6]. R. Yue *et al.* [7] applied a GCN that utilizes protein-protein interactions, drug-protein interactions, and drug-target interactions to predict potential drug molecules capable of binding with disease-related proteins.

Back-propagation (BP) was proposed by David Rumelhart, Geoffrey Hinton, and Ronald Williams as a learning procedure for training neural networks utilizing a set of input-output training samples [8, 9]. BP makes use of chain rule to compute the gradient of the network’s error with respect to every single model parameter, allowing for the adjustment of weight and bias terms via gradient descent to reduce the error [10, 11]. Reverse mode automatic differentiation is the primary technique in the form of the back-propagation algorithm for training neural networks due to its computational efficiency, particularly for an objective functions with a large number of inputs and a scalar output [12].

M. A. Nielsen [13] derived the BP for multi-layer perceptrons (MLPs) in a matrix-based form using the Hadamard product for numerical efficiency. N. M. Mishachev [14] utilized Hadamard product and Kronecker product to derive an explicit matrix version of the BP equations for MLPs and reduce the number of the indices in the equations. M. Naumov [15] demonstrated that the gradient of the weights can be expressed as a rank-1 and rank- t matrix for MLPs and recurrent neural networks at time step t , respectively. Y. Cheng [16] derived the BP algorithm for MLPs based on derivative amplification coefficients.

Explainable Artificial Intelligence (XAI) focuses on developing AI systems that not only make accurate predictions but also provide clear, interpretable explanations for their actions and decisions, thus increasing their trustworthiness for human users [17]. In XAI, sensitivity analysis is crucial because it reveals which parameters, or groups of parameters, have the most significant influence on the predictions made by machine learning models [18].

In this work, we used matrix calculus to derive an analytic and exact closed-form solution of the gradient of the loss function with respect to weight matrices for the training of GCN. The derivation considers the problems of binary node classification and link prediction. We first considered the problem of node classification with three-layer GCN and we extended the GCN model to a multi-layer GCN with d layers and applied the procedure for a link prediction problem. To validate our approach, we applied our weight gradient calculation for node classification on a Zachary’s Karate Club graph [19] and for link prediction using a self-prepared 100-node drug-drug interaction network. Subsequently, we compared the updated weight matrix using our method to the method employing reverse mode automatic differentiation for gradient computation. Our results demonstrate a ignorable difference between our weight matrix and the weight matrix updated using the gradient from reverse mode automatic differentiation, which verifies the correctness of our method. As there is no existing work deriving the BP of GCN in matrix form, we believe it is important for the machine learning community to have this closed-form solution available. In addition, we applied the same procedure to determine the sensitivity of the loss or output with respect to the input feature matrix in GCN for XAI. The definition of notations and some properties used in this paper are provided in Appendix A.

2 Back-propagation of Graph Convolutional Network

A basic graph structure is defined as:

$$G = (V, E), \quad (1)$$

where $|V| = n$ is the number of nodes in the graph and $|E| = n_e$ is the number of edges. Denoting $v_i \in V$ as a node and $e_{ij} = (v_i, v_j) \in E$ as an edge pointing from v_i to v_j , the adjacency matrix, $\mathbf{A} \in \{0, 1\}^{n \times n}$, is an $n \times n$ matrix where a_{ij} equals 1 if the edge e_{ij} exists, and a_{ij} equals 0 if e_{ij} does not belong to E ; in addition, a graph may possess node attributes represented by the matrix $\mathbf{H}_0 \in \mathbb{R}^{n \times n_0}$, where $\mathbf{h}_{0_v} \in \mathbb{R}^{n_0}$ is the feature vector of a node v with n_0 features [20].

2.1 Binary classification of nodes

2.1.1 Backpropagation for 3-layer GCN with ReLU and sigmoid activation function

We consider a 3-layer GCN defined as

$$\mathbf{H}_1 = \sigma_{ReLU}(\mathbf{A}\mathbf{H}_0\mathbf{W}_1), \quad (2)$$

$$\mathbf{H}_2 = \sigma_{ReLU}(\mathbf{A}\mathbf{H}_1\mathbf{W}_2), \quad (3)$$

$$\mathbf{H}_3 = \sigma_{ReLU}(\mathbf{A}\mathbf{H}_2\mathbf{W}_3), \quad (4)$$

$$\hat{\mathbf{Y}} = \sigma_{sigmoid}(\mathbf{H}_3), \quad (5)$$

where $\mathbf{A} \in \{0, 1\}^{n \times n}$ is an $n \times n$ adjacency matrix, $\mathbf{H}_0 \in \mathbb{R}^{n \times n_0}$ is the feature matrix for n nodes with n_0 features, $\hat{\mathbf{Y}} \in \mathbb{R}^{n \times n_3}$ is the output matrix, $n_3 = 1$ for binary node classification, and $\mathbf{W}_1 \in \mathbb{R}^{n_0 \times n_1}$, $\mathbf{W}_2 \in \mathbb{R}^{n_1 \times n_2}$, and $\mathbf{W}_3 \in \mathbb{R}^{n_2 \times n_3}$ are trainable parameter matrices, n_1 , n_2 , and n_3 are the number of features in $\mathbf{H}_1 \in \mathbb{R}^{n \times n_1}$, $\mathbf{H}_2 \in \mathbb{R}^{n \times n_2}$, and $\mathbf{H}_3 \in \mathbb{R}^{n \times n_3}$, respectively, σ_{ReLU} denotes an element-wise rectified linear unit (ReLU) function $\sigma_{ReLU}(v) = \max(v, 0)$ [21], $\sigma_{sigmoid}$ denotes an element-wise sigmoid function $\sigma_{sigmoid}(v) = \frac{1}{1+e^{-v}}$.

The loss function of the 3-layer GCN for binary node classification can be defined as

$$L = L_1 + L_2, \quad (6)$$

where

$$L_1 = - \sum_{i=1}^n \sum_{j=1}^{n_3} (y_{ij} \ln(\hat{y}_{ij})), \quad (7)$$

$$L_2 = - \sum_{i=1}^n \sum_{j=1}^{n_3} ((1 - y_{ij}) \ln(1 - \hat{y}_{ij})), \quad (8)$$

where y_{ij} denotes the element in the i th row and j th column of the ground truth of nodes $\mathbf{Y} \in \mathbb{R}^{n \times n_3}$, \hat{y}_{ij} denotes the element in the i th row and j th column of $\hat{\mathbf{Y}}$.

To derive the derivative of the loss function L with respect to the third weight matrix \mathbf{W}_3 , we first write

$$\frac{\partial L}{\partial \mathbf{W}_3} = \frac{\partial L_1}{\partial \mathbf{W}_3} + \frac{\partial L_2}{\partial \mathbf{W}_3}. \quad (9)$$

Then, the derivative of L_1 in (7) with respect to \mathbf{W}_3 can be written as

$$\begin{aligned} \frac{\partial L_1}{\partial \mathbf{W}_3} &= - \sum_{i=1}^n \sum_{j=1}^{n_3} y_{ij} \cdot \frac{\partial \ln(\hat{y}_{ij})}{\partial \mathbf{W}_3} \\ &= - \sum_{i=1}^n \sum_{j=1}^{n_3} y_{ij} \cdot \frac{\partial \ln(\sigma_{sigmoid}(h_{3ij}))}{\partial \mathbf{W}_3} \\ &= - \sum_{i=1}^n \sum_{j=1}^{n_3} y_{ij} \cdot \frac{1}{\hat{y}_{ij}} \cdot \hat{y}_{ij} \cdot (1 - \hat{y}_{ij}) \cdot \frac{\partial h_{3ij}}{\partial \mathbf{W}_3} \\ &= - \sum_{i=1}^n \sum_{j=1}^{n_3} y_{ij} \cdot (1 - \hat{y}_{ij}) \cdot \frac{\partial h_{3ij}}{\partial \mathbf{W}_3}. \end{aligned} \quad (10)$$

Similarly, the derivative of L_2 in (8) with respect to \mathbf{W}_3 can be written as

$$\begin{aligned} \frac{\partial L_2}{\partial \mathbf{W}_3} &= - \sum_{i=1}^n \sum_{j=1}^{n_3} (1 - y_{ij}) \cdot \frac{\partial \ln(1 - \hat{y}_{ij})}{\partial \mathbf{W}_3} \\ &= \sum_{i=1}^n \sum_{j=1}^{n_3} (1 - y_{ij}) \cdot \hat{y}_{ij} \cdot \frac{\partial h_{3ij}}{\partial \mathbf{W}_3}. \end{aligned} \quad (11)$$

The last term, $\frac{\partial h_{3ij}}{\partial \mathbf{W}_3}$, in (10) and (11) can be written as $\frac{\partial h_{3ij}}{\partial \mathbf{W}_3} = \frac{\partial \sigma_{ReLU}(\mathbf{A}_{i*}\mathbf{H}_2\mathbf{W}_{3*j})}{\partial \mathbf{W}_3}$ using (4), where \mathbf{A}_{i*} is the i th row of \mathbf{A} and \mathbf{W}_{3*j} is the j th column of \mathbf{W}_3 . The expression, $\frac{\partial \sigma_{ReLU}(\mathbf{A}_{i*}\mathbf{H}_2\mathbf{W}_{3*j})}{\partial \mathbf{W}_3}$, is solved using the following definitions and theorem.

Definition 2.1. An element-wise activation function $\Sigma : \mathbb{R}^{p \times q} \rightarrow \mathbb{R}^{p \times q}$ is defined as a multivariate matrix-valued function [22]:

$$\Sigma : \mathbf{X} \mapsto \Sigma(\mathbf{X}) := \begin{bmatrix} \sigma(x_{11}) & \cdots & \sigma(x_{1q}) \\ \vdots & & \vdots \\ \sigma(x_{p1}) & \cdots & \sigma(x_{pq}) \end{bmatrix}, \quad (12)$$

where $\mathbf{X} \in \mathbb{R}^{p \times q}$ and $\sigma : \mathbb{R} \rightarrow \mathbb{R}$, $x_{ij} \mapsto \sigma(x_{ij})$ for all $i = 1, 2, \dots, p$ and $j = 1, 2, \dots, q$.

Definition 2.2. The derivative of an element-wise activation function $\Sigma(\mathbf{X})$ ($p \times q$) with respect to a matrix $\mathbf{X} \in \mathbb{R}^{p \times q}$ is written as:

$$\Sigma'(\mathbf{X}) = \frac{\partial \Sigma(\mathbf{X})}{\partial \mathbf{X}} = \begin{bmatrix} \sigma'(x_{11}) & \cdots & \sigma'(x_{1q}) \\ \vdots & & \vdots \\ \sigma'(x_{p1}) & \cdots & \sigma'(x_{pq}) \end{bmatrix}, \quad (13)$$

where $\sigma'(v) = \lim_{u \rightarrow 0} \frac{\sigma(v+u) - \sigma(v)}{u}$ is defined as the derivative of a real-valued function σ .

Theorem 2.3. Let $\mathbf{F} : \mathbb{R}^{p \times q} \rightarrow \mathbb{R}^{m \times n}$ be a $m \times n$ multivariate matrix-valued function of a $p \times q$ matrix $\mathbf{W} \in \mathbb{R}^{p \times q}$, the derivative of $\Sigma(\mathbf{F}(\mathbf{W}))$ with respect to \mathbf{W} can be written as:

$$\frac{\partial \Sigma(\mathbf{F}(\mathbf{W}))}{\partial \mathbf{W}} = (\mathbf{J}_{p \times q} \otimes \Sigma'(\mathbf{F}(\mathbf{W}))) \odot \frac{\partial \mathbf{F}(\mathbf{W})}{\partial \mathbf{W}}, \quad (14)$$

where \otimes is the Kronecker product in (46), \odot is the Hadamard product in (47), and $\mathbf{J}_{p \times q} \in \mathbb{R}^{p \times q}$ is an all 1 matrix.

Proof. The proof can be found in Appendix C. □

From theorem 2.3, we can write the last term in (10) and (11) as

$$\begin{aligned} \frac{\partial h_{3ij}}{\partial \mathbf{W}_3} &= \frac{\partial \sigma_{ReLU}(\mathbf{A}_{i*} \mathbf{H}_2 \mathbf{W}_{3*j})}{\partial \mathbf{W}_3} \\ &= (\mathbf{J}_{n_2 \times n_3} \otimes \sigma'_{ReLU}(\mathbf{A}_{i*} \mathbf{H}_2 \mathbf{W}_{3*j})) \odot \frac{\partial \mathbf{A}_{i*} \mathbf{H}_2 \mathbf{W}_{3*j}}{\partial \mathbf{W}_3}. \end{aligned} \quad (15)$$

Using (52), the last term in (15) can be written as

$$\begin{aligned} \frac{\partial \mathbf{A}_{i*} \mathbf{H}_2 \mathbf{W}_{3*j}}{\partial \mathbf{W}_3} &= \frac{\partial \mathbf{A}_{i*} \mathbf{H}_2}{\partial \mathbf{W}_3} (\mathbf{I}_{n_3} \otimes \mathbf{W}_{3*j}) + (\mathbf{I}_{n_2} \otimes (\mathbf{A}_{i*} \mathbf{H}_2)) \frac{\partial \mathbf{W}_{3*j}}{\partial \mathbf{W}_3} \\ &= (\mathbf{I}_{n_2} \otimes (\mathbf{A}_{i*} \mathbf{H}_2)) \frac{\partial \mathbf{W}_{3*j}}{\partial \mathbf{W}_3}. \end{aligned} \quad (16)$$

Using (44), (52), and (53), the last term in (16) can be written as

$$\begin{aligned} \frac{\partial \mathbf{W}_{3*j}}{\partial \mathbf{W}_3} &= \frac{\partial \mathbf{W}_3 \mathbf{e}_j}{\partial \mathbf{W}_3} \\ &= \frac{\partial \mathbf{W}_3}{\partial \mathbf{W}_3} (\mathbf{I}_{n_3} \otimes \mathbf{e}_j) + (\mathbf{I}_{n_2} \otimes \mathbf{W}_3) \frac{\partial \mathbf{e}_j}{\partial \mathbf{W}_3} \\ &= \bar{\mathbf{U}}_{n_2 \times n_3} (\mathbf{I}_{n_3} \otimes \mathbf{e}_j), \end{aligned} \quad (17)$$

where \mathbf{e}_j is a n_3 -dimensional column vector which has “1” in the j th row and zero elsewhere and

$\bar{\mathbf{U}}_{n_2 \times n_3} \in \mathbb{R}^{n_2^2 \times n_3^2}$ is a permutation related matrix defined in (49).

Thus, the derivative of the loss function L with respect to the third weight matrix \mathbf{W}_3 in (9) can be written as

$$\frac{\partial L}{\partial \mathbf{W}_3} = - \sum_{i=1}^n \sum_{j=1}^{n_3} ((y_{ij} - \hat{y}_{ij}) \cdot \frac{\partial h_{3ij}}{\partial \mathbf{W}_3}) \quad (18)$$

$$\begin{aligned} &= - \sum_{i=1}^n \sum_{j=1}^{n_3} ((y_{ij} - \hat{y}_{ij}) \cdot ((\mathbf{J}_{n_2 \times n_3} \otimes \sigma'_{ReLU}(\mathbf{A}_{i*} \mathbf{H}_2 \mathbf{W}_{3*j})) \\ &\quad \odot ((\mathbf{I}_{n_2} \otimes (\mathbf{A}_{i*} \mathbf{H}_2)) \bar{\mathbf{U}}_{n_2 \times n_3} (\mathbf{I}_{n_3} \otimes \mathbf{e}_j))))). \end{aligned} \quad (19)$$

Using (18), the derivative of the loss function L in (6) with respect to the second weight matrix \mathbf{W}_2 can be written as

$$\frac{\partial L}{\partial \mathbf{W}_2} = - \sum_{i=1}^n \sum_{j=1}^{n_3} ((y_{ij} - \hat{y}_{ij}) \cdot \frac{\partial h_{3ij}}{\partial \mathbf{W}_2}) \quad (20)$$

$$= - \sum_{i=1}^n \sum_{j=1}^{n_3} ((y_{ij} - \hat{y}_{ij}) \cdot ((\mathbf{J}_{n_1 \times n_2} \otimes \sigma'_{ReLU}(\mathbf{A}_{i*} \mathbf{H}_2 \mathbf{W}_{3*j})) \odot \frac{\partial \mathbf{A}_{i*} \mathbf{H}_2 \mathbf{W}_{3*j}}{\partial \mathbf{W}_2}))), \quad (21)$$

where (21) follows from theorem 2.3.

Using (52), the last term of (21) can be written as

$$\frac{\partial \mathbf{A}_{i*} \mathbf{H}_2 \mathbf{W}_{3*j}}{\partial \mathbf{W}_2} = (\mathbf{I}_{n_1} \otimes \mathbf{A}_{i*}) \frac{\partial \mathbf{H}_2}{\partial \mathbf{W}_2} (\mathbf{I}_{n_2} \otimes \mathbf{W}_{3*j}). \quad (22)$$

The derivative of \mathbf{H}_2 with respect to \mathbf{W}_2 in (22) can be written as

$$\frac{\partial \mathbf{H}_2}{\partial \mathbf{W}_2} = \frac{\partial \sigma_{ReLU}(\mathbf{A} \mathbf{H}_1 \mathbf{W}_2)}{\partial \mathbf{W}_2} \quad (23)$$

$$= (\mathbf{J}_{n_1 \times n_2} \otimes \sigma'_{ReLU}(\mathbf{A} \mathbf{H}_1 \mathbf{W}_2)) \odot \frac{\partial \mathbf{A} \mathbf{H}_1 \mathbf{W}_2}{\partial \mathbf{W}_2} \quad (24)$$

$$= (\mathbf{J}_{n_1 \times n_2} \otimes \sigma'_{ReLU}(\mathbf{A} \mathbf{H}_1 \mathbf{W}_2)) \odot ((\mathbf{I}_{n_1} \otimes (\mathbf{A} \mathbf{H}_1)) \frac{\partial \mathbf{W}_2}{\partial \mathbf{W}_2}) \quad (25)$$

$$= (\mathbf{J}_{n_1 \times n_2} \otimes \sigma'_{ReLU}(\mathbf{A} \mathbf{H}_1 \mathbf{W}_2)) \odot ((\mathbf{I}_{n_1} \otimes (\mathbf{A} \mathbf{H}_1)) \bar{\mathbf{U}}_{n_1 \times n_2}), \quad (26)$$

where (23) follows from (4), (24) follows from (15), (25) follows from (52), and (26) follows from (53).

The derivative of the loss in (6) with respect to the first weight matrix \mathbf{W}_1 is

$$\begin{aligned} \frac{\partial L}{\partial \mathbf{W}_1} &= - \sum_{i=1}^n \sum_{j=1}^{n_3} ((y_{ij} - \hat{y}_{ij}) \cdot ((\mathbf{J}_{n_0 \times n_1} \otimes \sigma'_{ReLU}(\mathbf{A}_{i*} \mathbf{H}_2 \mathbf{W}_{3*j})) \\ &\quad \odot ((\mathbf{I}_{n_0} \otimes \mathbf{A}_{i*}) \frac{\partial \mathbf{H}_2}{\partial \mathbf{W}_1} (\mathbf{I}_{n_1} \otimes \mathbf{W}_{3*j}))))), \end{aligned} \quad (27)$$

where (27) follows from (21) and

$$\frac{\partial \mathbf{H}_2}{\partial \mathbf{W}_1} = \frac{\partial \sigma_{ReLU}(\mathbf{A} \mathbf{H}_1 \mathbf{W}_2)}{\partial \mathbf{W}_1} \quad (28)$$

$$= (\mathbf{J}_{n_0 \times n_1} \otimes \sigma'_{ReLU}(\mathbf{A} \mathbf{H}_1 \mathbf{W}_2)) \odot ((\mathbf{I}_{n_0} \otimes \mathbf{A}) \frac{\partial \mathbf{H}_1}{\partial \mathbf{W}_1} (\mathbf{I}_{n_1} \otimes \mathbf{W}_2)), \quad (29)$$

where (28) follows from (4), (29) follows from theorem (2.3) and (52), and

$$\frac{\partial \mathbf{H}_1}{\partial \mathbf{W}_1} = (\mathbf{J}_{n_0 \times n_1} \otimes \sigma'_{ReLU}(\mathbf{A} \mathbf{H}_0 \mathbf{W}_1)) \odot ((\mathbf{I}_{n_0} \otimes (\mathbf{A} \mathbf{H}_0)) \bar{\mathbf{U}}_{n_0 \times n_1}), \quad (30)$$

where (30) follows from (26).

2.1.2 Back-propagation for multi-layer GCN with ReLU and sigmoid activation function

We consider a d -layer GCN defined as

$$\hat{\mathbf{Y}} = \sigma_{\text{sigmoid}}(\mathbf{H}_d), \quad (31)$$

where

$$\mathbf{H}_d = \sigma_{\text{ReLU}}(\mathbf{A}\mathbf{H}_{d-1}\mathbf{W}_d), \quad (32)$$

$d \in \mathbb{Z}^+$, $\mathbf{A} \in \{0, 1\}^{n \times n}$ is an $n \times n$ adjacency matrix, $\mathbf{H}_{d-1} \in \mathbb{R}^{n \times n_{d-1}}$ is the feature matrix for n nodes with n_{d-1} features, $\hat{\mathbf{Y}} \in \mathbb{R}^{n \times n_d}$ is the output matrix, $n_d = 1$ for binary node classification, and $\mathbf{W}_1 \in \mathbb{R}^{n_0 \times n_1}$, $\mathbf{W}_2 \in \mathbb{R}^{n_1 \times n_2}$, \dots , and $\mathbf{W}_d \in \mathbb{R}^{n_{d-1} \times n_d}$ are trainable parameter matrices.

By observing (19), (21), (22), (26), (27), (29), and (30), if the loss function of the d -layer GCN for node classification is defined as

$$L = - \sum_{i=1}^n \sum_{j=1}^{n_d} (y_{ij} \ln(\hat{y}_{ij}) + (1 - y_{ij}) \ln(1 - \hat{y}_{ij})), \quad (33)$$

the derivative of the loss in (33) with respect to the s th weight matrix \mathbf{W}_s is

$$\begin{aligned} \frac{\partial L}{\partial \mathbf{W}_s} = & - \sum_{i=1}^n \sum_{j=1}^{n_d} ((y_{ij} - \hat{y}_{ij}) \cdot ((\mathbf{J}_{(s-1) \times s} \otimes \sigma'_{\text{ReLU}}(\mathbf{A}_{i*} \mathbf{H}_{d-1} \mathbf{W}_{d*j})) \\ & \odot \frac{\partial \mathbf{A}_{i*} \mathbf{H}_{d-1} \mathbf{W}_{d*j}}{\partial \mathbf{W}_s})), \end{aligned} \quad (34)$$

where $s \in \mathbb{Z}^+$, $s \leq d$,

$$\frac{\partial \mathbf{A}_{i*} \mathbf{H}_{d-1} \mathbf{W}_{d*j}}{\partial \mathbf{W}_s} = \begin{cases} (\mathbf{I}_{n_{s-1}} \otimes (\mathbf{A}_{i*} \mathbf{H}_{d-1})) \cdot \bar{\mathbf{U}}_{n_{s-1} \times n_s} (\mathbf{I}_{n_s} \otimes \underset{(n_d)}{\mathbf{e}_j}), & \text{if } s = d \\ (\mathbf{I}_{n_{s-1}} \otimes \mathbf{A}_{i*}) \cdot \frac{\partial \mathbf{H}_{d-1}}{\partial \mathbf{W}_s} (\mathbf{I}_{n_s} \otimes \mathbf{W}_{d*j}), & \text{if } s < d \end{cases} \quad (35)$$

and

$$\begin{aligned} & \frac{\partial \mathbf{H}_{d-1}}{\partial \mathbf{W}_s} \\ = & \begin{cases} (\mathbf{J}_{n_{s-1} \times n_s} \otimes \sigma'_{\text{ReLU}}(\mathbf{A}\mathbf{H}_{d-2}\mathbf{W}_{d-1})) \odot ((\mathbf{I}_{n_{s-1}} \otimes (\mathbf{A}\mathbf{H}_{d-2})) \bar{\mathbf{U}}_{n_{s-1} \times n_s}), & \text{if } s = d-1 \\ (\mathbf{J}_{n_{s-1} \times n_s} \otimes \sigma'_{\text{ReLU}}(\mathbf{A}\mathbf{H}_{d-2}\mathbf{W}_{d-1})) \odot ((\mathbf{I}_{n_{s-1}} \otimes \mathbf{A}) \frac{\partial \mathbf{H}_{d-2}}{\partial \mathbf{W}_s} (\mathbf{I}_{n_s} \otimes \mathbf{W}_{d-1})), & \text{if } s < d-1. \end{cases} \end{aligned} \quad (36)$$

2.1.3 Back-propagation for multi-layer GCN with arbitrary activation functions

We consider a d -layer GCN defined as

$$\hat{\mathbf{Y}} = \Sigma_{d+1}(\mathbf{H}_d), \quad (37)$$

where

$$\mathbf{H}_d = \Sigma_d(\mathbf{A}\mathbf{H}_{d-1}\mathbf{W}_d), \quad (38)$$

$d \in \mathbb{Z}^+$, $\mathbf{A} \in \{0, 1\}^{n \times n}$ is an $n \times n$ adjacency matrix, $\mathbf{H}_{d-1} \in \mathbb{R}^{n \times n_{d-1}}$ is the feature matrix for n nodes with n_{d-1} features, $\hat{\mathbf{Y}} \in \mathbb{R}^{n \times n_d}$ is the output matrix, $n_d = 1$ for binary node classification, $\mathbf{W}_1 \in \mathbb{R}^{n_0 \times n_1}$, $\mathbf{W}_2 \in \mathbb{R}^{n_1 \times n_2}$, \dots , and $\mathbf{W}_d \in \mathbb{R}^{n_{d-1} \times n_d}$ are trainable parameter matrices, and $\Sigma_1, \Sigma_2, \dots$, and Σ_{d+1} are any element-wise activation function.

If the loss function of the d -layer GCN for node classification is defined as

$$L = - \sum_{i=1}^n \sum_{j=1}^{n_d} (y_{ij} \ln(\hat{y}_{ij}) + (1 - y_{ij}) \ln(1 - \hat{y}_{ij})), \quad (39)$$

the derivative of the loss in (39) with respect to the s th weight matrix \mathbf{W}_s is

$$\begin{aligned}
\frac{\partial L}{\partial \mathbf{W}_s} &= - \sum_{i=1}^n \sum_{j=1}^{n_d} (y_{ij} \frac{\partial \ln(\hat{y}_{ij})}{\partial \mathbf{W}_s} + (1 - \hat{y}_{ij}) \frac{\partial \ln(1 - \hat{y}_{ij})}{\partial \mathbf{W}_s}) \\
&= - \sum_{i=1}^n \sum_{j=1}^{n_d} ((\frac{y_{ij}}{\hat{y}_{ij}} - \frac{1 - y_{ij}}{1 - \hat{y}_{ij}}) \Sigma'_{d+1}(h_{d_{ij}}) \cdot \frac{\partial \Sigma_d(\mathbf{A}_{i*} \mathbf{H}_{d-1} \mathbf{W}_{d*j})}{\partial \mathbf{W}_s}) \\
&= - \sum_{i=1}^n \sum_{j=1}^{n_d} (\frac{y_{ij} - \hat{y}_{ij}}{\hat{y}_{ij}(1 - \hat{y}_{ij})} \Sigma'_{d+1}(h_{d_{ij}}) \cdot ((\mathbf{J}_{(s-1) \times s} \otimes \Sigma'_d(\mathbf{A}_{i*} \mathbf{H}_{d-1} \mathbf{W}_{d*j})) \\
&\quad \odot \frac{\partial \mathbf{A}_{i*} \mathbf{H}_{d-1} \mathbf{W}_{d*j}}{\partial \mathbf{W}_s})),
\end{aligned} \tag{40}$$

$$\tag{41}$$

where $s \in \mathbb{Z}^+$, $s \leq d$, (40) follows from chain rule and definition 2.2, (41) follows from theorem (2.3), $\frac{\partial \mathbf{A}_{i*} \mathbf{H}_{d-1} \mathbf{W}_{d*j}}{\partial \mathbf{W}_s}$ is the same as in (35), and

$$\begin{aligned}
&\frac{\partial \mathbf{H}_{d-1}}{\partial \mathbf{W}_s} \\
&= \begin{cases} (\mathbf{J}_{n_{s-1} \times n_s} \otimes \Sigma'_{d-1}(\mathbf{A} \mathbf{H}_{d-2} \mathbf{W}_{d-1})) \odot ((\mathbf{I}_{n_{s-1}} \otimes (\mathbf{A} \mathbf{H}_{d-2})) \bar{\mathbf{U}}_{n_{s-1} \times n_s}) & \text{if } s = d - 1 \\ (\mathbf{J}_{n_{s-1} \times n_s} \otimes \Sigma'_{d-1}(\mathbf{A} \mathbf{H}_{d-2} \mathbf{W}_{d-1})) \odot ((\mathbf{I}_{n_{s-1}} \otimes \mathbf{A}) \frac{\partial \mathbf{H}_{d-2}}{\partial \mathbf{W}_s} (\mathbf{I}_{n_s} \otimes \mathbf{W}_{d-1})) & \text{if } s < d - 1. \end{cases}
\end{aligned} \tag{42}$$

Following a similar procedure, we can derive the sensitivity of the loss with respect to the feature matrix \mathbf{H}_0 , as shown in Appendix E.1.

2.2 Link prediction

For the derivation of back-propagation and sensitivity in a multi-layer GCN with arbitrary activation functions for the link prediction problem, please refer to Appendices D and E.2, respectively.

3 Experiments

In this section, we demonstrate the correctness of our method through several experiments considering a node classification problem on Zachary's Karate Club [19] in Figure 3 (a) and a link prediction problem on a drug-drug interaction (DDI) network in Figure 3 (b). We adopt the stochastic gradient descent (SGD) for model training. We compare our method to the gradient computation using reverse mode automatic differentiation in Pytorch [23] by computing the sum of squared error (SSE) defined as

$$SSE = \sum_{i=1}^{n_{s-1}} \sum_{j=1}^{n_s} (w_{s_{ij}}^{(AD)} - w_{s_{ij}}^{(KP)})^2, \tag{43}$$

where $w_{s_{ij}}^{(AD)}$ is the i th row and the j th column of the s th weight matrix \mathbf{W}_s updated using SGD with reverse mode automatic differentiation and $w_{s_{ij}}^{(KP)}$ is the i th row and the j th column of the s th weight matrix \mathbf{W}_s updated using SGD with our matrix-based method.

In addition, we determine the sensitivity of the loss with respect to the feature matrix \mathbf{H}_0 for node classification and link prediction to demonstrate the application of our method to XAI.

The python codes are available at: <https://github.com/AnonymousForPapers/GCN-proof/tree/main>.

3.1 Node classification

We tested our method for node classification on Zachary's Karate Club [19] (see Appendix F.1 for more details).

3.1.1 1-layer GCN with identity function and sigmoid activation function

The 1-layer GCN is defined by (37) and the loss function is defined by (39), where $d = 1$ is the number of layers, $\mathbf{H}_0 = \mathbf{I}_n \in \mathbb{R}^{n \times n}$ is the feature matrix for $n = 34$ nodes with $n = 34$ features, the matrix $\hat{\mathbf{Y}} \in \mathbb{R}^{n \times n_d}$ represents the ground truth of nodes, indicating class assignments per node, $\hat{\mathbf{Y}} \in \mathbb{R}^{n \times n_d}$ denotes the final layer predictions of the GCN model, $n_d = 1$ for binary node classification, $\mathbf{W}_1 \in \mathbb{R}^{n_0 \times n_d}$ is a trainable parameter matrix, $n_0 = 34$, Σ_1 is an element-wise identity function, and Σ_2 is an element-wise sigmoid function.

The GCN is trained using SGD with a learning rate of 0.1 and with 100 iterations for both methods. In Figure 4, the evolution of the karate club network shows that the two graph updated by the two methods have the same loss, accuracy, and classification results. In Figure 1 (a), it shows the SSE of the weight matrix \mathbf{W}_1 between the reverse mode automatic differentiation and our matrix-based method is lower than 10^{-13} . This indicates that our analytic expression aligns with the exact method, and the small difference might be due to the precision of different datatypes.

Next, we show the loss sensitivity with respect to the input feature matrix, computed using the result in Appendix E.1. In Figure 2 (a), it is shown that as training progresses, the sensitivity decreases. This result is expected since the input feature matrix of the GCN for the Karate graph is an identity matrix. The prediction of the class of a node depends only on its edges rather than the input features. Therefore, changes in the input matrix should not significantly affect the prediction from the trained GCN, as it should learn to make predictions based solely on the existence of edges between nodes.

Additional experimental results on Zachary’s Karate Club [19] are provided in Appendix G.0.1.

3.2 Link prediction

We tested our method for link prediction on a 10-node DDI network (see Appendix F.2 for more details).

3.2.1 2-layer GCN

The 2-layer GCN is defined by (60) and the loss function is defined by (62), where $d = 2$ is the number of layers, the matrix $\hat{\mathbf{Y}} \in \mathbb{R}^{n \times n}$ represents ground-truth adjacency matrix, $\hat{\mathbf{Y}} \in \mathbb{R}^{n \times n}$ denotes the final link predictions of the GCN model, $\mathbf{W}_1 \in \mathbb{R}^{n_0 \times n_1}$ and $\mathbf{W}_2 \in \mathbb{R}^{n_1 \times n_2}$ are trainable parameter matrices, $n_0 = 20$, $n_1 = 10$, $n_2 = 5$, Σ_1 is an element-wise ReLU function [21], Σ_2 is an element-wise identity function, and Σ_3 is an element-wise sigmoid function.

The GCN is trained using SGD with a learning rate of 0.01 and with 150 iterations for both methods. In each iteration, thirteen negative edges, which is the same number as the number of positive edges, are uniformly sampled from a set of all the unconnected edges for the GCN training using the reverse mode automatic differentiation. The sampled negative edges in each iteration are saved in a list and used in the training of the GCN using our matrix-based method. In Figure 7, the evolution of the DDI network shows that the two graph updated by the two methods have the same loss and classification results. In Figure 1 (b), it shows the SSE of the two weight matrices, \mathbf{W}_1 and \mathbf{W}_2 , between the reverse mode automatic differentiation and our matrix-based method is lower than 10^{-14} .

Next, we show the output sensitivity with respect to the input feature matrix, computed using the result in Appendix E.2. In Figure 2 (b), a heat map of the sensitivity of the prediction for the link between node 2 and node 7 is shown. The sensitivities of the features in node 2 and node 7 are higher than those of the features in other nodes. The features of nodes (node 1 and node 8) that do not have a connection to either node 2 or node 7, or the features of nodes (node 0) that need to traverse at least three edges to reach node 2 or node 7, have almost zero sensitivity to the link. This demonstrates the property of GCNs, where the neighborhood information of two nodes is aggregated by taking the weighted sum of the features of neighboring nodes [24]. Since the GCN has only two layers, the information of nodes that are three edges away from these two nodes does not contribute to the prediction of the edge.

Additional experimental results on the DDI network are provided in Appendix G.0.2. Heat maps of the sensitivity of the prediction for all the links are in Figures 8, 9, and 10 in the Appendix.

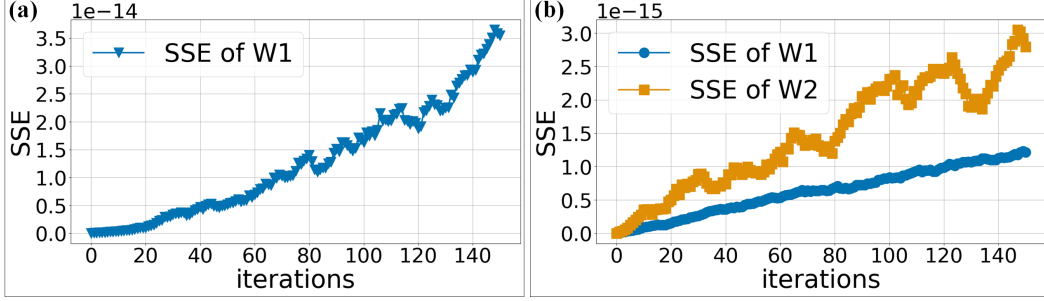


Figure 1: (a) The evolution of the sum of squared error between the trainable weight matrix obtained from our method and the matrix obtained using reverse mode automatic differentiation in section 3.1.1. (b) The evolution of the sum of squared error between the two trainable weight matrices obtained from our method and the matrices obtained using reverse mode automatic differentiation in section 3.2.1.

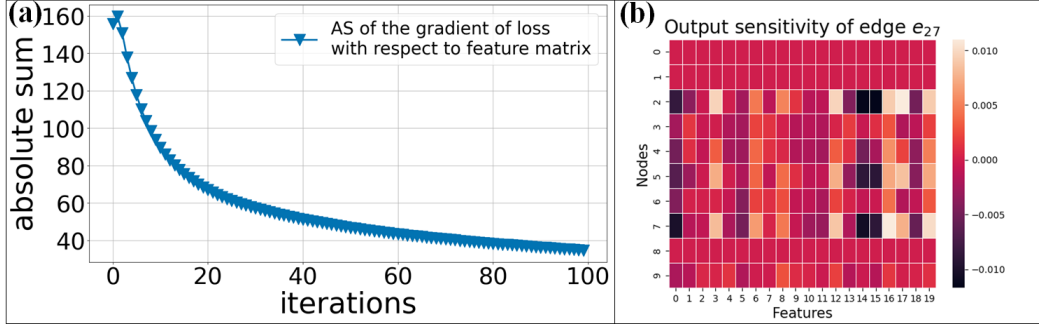


Figure 2: (a) The evolution of the absolute sum of the sensitivity of the loss with respect to the input feature matrix \mathbf{H}_0 in Section 3.1.1. (b) The heat map of the sensitivity of the prediction for the link between node 2 and node 7 with respect to the input feature matrix \mathbf{H}_0 in Section 3.2.1.

4 Conclusion

In this work, we provide detailed derivations of the analytical expressions in matrix form for the derivatives of a loss function with respect to each weight matrix for a graph convolutional network considering binary node classification and link prediction. Utilizing Kronecker product, Hadamard product, and matrix calculus, we computed the gradients of the loss function in a three-layer graph convolutional network and extended the solution to accommodate graph convolutional networks with arbitrary layers and element-wise activation functions. The weight matrices obtained through our approach were compared with those obtained using reverse mode automatic differentiation in binary node classification experiments conducted on a 34 nodes Zachary’s Karate Club network [19] and in link prediction experiments using a 10-node drug-drug interaction network. The experimental results indicate that the discrepancy between the weight matrices has a median sum of squared error ranging from 10^{-18} to 10^{-14} , affirming the accuracy of our methodology. In addition, we conduct sensitivity analysis for binary node classification and link prediction to demonstrate the application of our derivation method in XAI.

We note that our derivation already accommodates RNNs when \mathbf{A} is an identity matrix, and it accommodates CNNs since CNNs operate on 2-dimensional matrices, which are a special case of graphs. Additionally, our method incurs a significantly higher computational cost compared to reverse mode AD. As future work, we aim to find a way to improve the computational speed of back-propagation through our derived analytical solution.

References

- [1] Franco Scarselli, Marco Gori, Ah Chung Tsoi, Markus Hagenbuchner, and Gabriele Monfardini. The graph neural network model. *IEEE transactions on neural networks*, 20(1):61–80, 2008.
- [2] Thomas N Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. *arXiv preprint arXiv:1609.02907*, 2016.
- [3] Abhishek Dutta. Deep graph generation of small molecules guided by drug-likeness. In *IEEE Conference on Computational Intelligence in Bioinformatics and Computational Biology*, pages 1–2. IEEE, 2022.
- [4] Shunxin Xiao, Shiping Wang, Yuanfei Dai, and Wenzhong Guo. Graph neural networks in node classification: survey and evaluation. *Machine Vision and Applications*, 33:1–19, 2022.
- [5] Jinxian Zhang, Junbo Zhao, Fei Ding, Jing Yang, and Junhui Zhao. A graph convolutional network for active distribution system anomaly detection considering measurement spatial-temporal correlations. In *2023 North American Power Symposium (NAPS)*, pages 1–6. IEEE, 2023.
- [6] Xing Wang and Alexander Vinel. Benchmarking graph neural networks on link prediction. *arXiv preprint arXiv:2102.12557*, 2021.
- [7] Rongting Yue and Abhishek Dutta. Repurposing drugs for covid-19 by graph convolutional network. *PREPRINT*, April 12 2023. Version 1, available at Research Square <https://doi.org/10.21203/rs.3.rs-2408594/v1>.
- [8] David E Rumelhart, Geoffrey E Hinton, and Ronald J Williams. Learning representations by back-propagating errors. *nature*, 323(6088):533–536, 1986.
- [9] Sergios Theodoridis. Chapter 18 - neural networks and deep learning. In Sergios Theodoridis, editor, *Machine Learning (Second Edition)*, pages 901–1038. Academic Press, second edition, 2020.
- [10] Aurélien Géron. *Hands-on machine learning with Scikit-Learn, Keras, and TensorFlow*. "O'Reilly Media, Inc.", 2022.
- [11] Werbos. Backpropagation: Past and future. In *IEEE 1988 International Conference on Neural Networks*, pages 343–353. IEEE, 1988.
- [12] Atilim Gunes Baydin, Barak A Pearlmutter, Alexey Andreyevich Radul, and Jeffrey Mark Siskind. Automatic differentiation in machine learning: a survey. *Journal of Machine Learning Research*, 18:1–43, 2018.
- [13] Michael A Nielsen. *Neural networks and deep learning*, volume 25. Determination press San Francisco, CA, USA, 2015.
- [14] NM Mishachev. Backpropagation in matrix notation. *arXiv preprint arXiv:1707.02746*, 2017.
- [15] Maxim Naumov. Feedforward and recurrent neural networks backward propagation and hessian in matrix form. *arXiv preprint arXiv:1709.06080*, 2017.
- [16] Yiping Cheng. Derivation of the backpropagation algorithm based on derivative amplification coefficients. *arXiv preprint arXiv:2102.04320*, 2021.
- [17] Vikas Hassija, Vinay Chamola, Atmesh Mahapatra, Abhinandan Singal, Divyansh Goel, Kaizhu Huang, Simone Scardapane, Indro Spinelli, Mufti Mahmud, and Amir Hussain. Interpreting black-box models: a review on explainable artificial intelligence. *Cognitive Computation*, 16(1):45–74, 2024.
- [18] Bas Van Stein, Elena Raponi, Zahra Sadeghi, Niek Bouman, Roeland CHJ Van Ham, and Thomas Bäck. A comparison of global sensitivity analysis methods for explainable ai with an application in genomic prediction. *IEEE Access*, 10:103364–103381, 2022.
- [19] Wayne W Zachary. An information flow model for conflict and fission in small groups. *Journal of anthropological research*, 33(4):452–473, 1977.
- [20] Zonghan Wu, Shirui Pan, Fengwen Chen, Guodong Long, Chengqi Zhang, and S Yu Philip. A comprehensive survey on graph neural networks. *IEEE transactions on neural networks and learning systems*, 32(1):4–24, 2020.

- [21] Vinod Nair and Geoffrey E Hinton. Rectified linear units improve restricted boltzmann machines. In *Proceedings of the 27th international conference on machine learning (ICML-10)*, pages 807–814, 2010.
- [22] Dirk Ostwald and Franziska Usée. An induction proof of the backpropagation algorithm in matrix notation. *arXiv preprint arXiv:2107.09384*, 2021.
- [23] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems*, 32, 2019.
- [24] David Ahmedt-Aristizabal, Mohammad Ali Armin, Simon Denman, Clinton Fookes, and Lars Petersson. A survey on graph-based deep learning for computational histopathology. *Computerized Medical Imaging and Graphics*, 95:102027, 2022.
- [25] John Brewer. Kronecker products and matrix calculus in system theory. *IEEE Transactions on circuits and systems*, 25(9):772–781, 1978.
- [26] Hai-hua Zhu, Zi-gang Chen, and Tao Leng. Random permutation-based mixed-double scrambling technique for encrypting mqir image. *Journal of Applied Physics*, 135(1), 2024.
- [27] Jan R Magnus. On the concept of matrix derivative. *Journal of Multivariate Analysis*, 101(9):2200–2206, 2010.
- [28] Seongjun Yun, Seoyoon Kim, Junhyun Lee, Jaewoo Kang, and Hyunwoo J Kim. Neo-gnns: Neighborhood overlap-aware graph neural networks for link prediction. *Advances in Neural Information Processing Systems*, 34:13683–13694, 2021.
- [29] Cristopher Moore, Xiaoran Yan, Yaojia Zhu, Jean-Baptiste Rouquier, and Terran Lane. Active learning for node classification in assortative and disassortative networks. In *Proceedings of the 17th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 841–849, 2011.
- [30] A David Rodrigues. *Drug-drug interactions*. CRC Press, 2019.
- [31] David S Wishart, Yannick D Feunang, An C Guo, Elvis J Lo, Ana Marcu, Jason R Grant, Tanvir Sajed, Daniel Johnson, Carin Li, Zinat Sayeeda, et al. Drugbank 5.0: a major update to the drugbank database for 2018. *Nucleic acids research*, 46(D1):D1074–D1082, 2018.
- [32] Daniel Blanco-Melo, Benjamin E Nilsson-Payant, Wen-Chun Liu, Skyler Uhl, Daisy Hoagland, Rasmus Møller, Tristan X Jordan, Kohei Oishi, Maryline Panis, David Sachs, et al. Imbalanced host response to sars-cov-2 drives development of covid-19. *Cell*, 181(5):1036–1045, 2020.
- [33] Minoru Kanehisa, Yoko Sato, and Masayuki Kawashima. Kegg mapping tools for uncovering hidden features in biological data. *Protein Science*, 2021.
- [34] David Weininger. Smiles, a chemical language and information system. 1. introduction to methodology and encoding rules. *Journal of chemical information and computer sciences*, 28(1):31–36, 1988.
- [35] Dan Hendrycks and Kevin Gimpel. Gaussian error linear units (gelus). *arXiv preprint arXiv:1606.08415*, 2016.
- [36] Djork-Arné Clevert, Thomas Unterthiner, and Sepp Hochreiter. Fast and accurate deep network learning by exponential linear units (elus). *arXiv preprint arXiv:1511.07289*, 2015.
- [37] Andrew L Maas, Awni Y Hannun, Andrew Y Ng, et al. Rectifier nonlinearities improve neural network acoustic models. In *Proc. icml*, volume 30, page 3. Atlanta, GA, 2013.
- [38] J. D. Hunter. Matplotlib: A 2d graphics environment. *Computing in Science & Engineering*, 9(3):90–95, 2007.

A Basic notation and properties of Kronecker product and matrix calculus

A.1 Basic notation

A column vector is denoted by lower case boldface (e.g., \mathbf{v} , with its i th element being \mathbf{v}_i). Matrices is denoted by upper case boldface (e.g., \mathbf{A}). The i th row for a matrix such as \mathbf{A} is denoted \mathbf{A}_{i*} and the i th column is denoted \mathbf{A}_{*i} . The (i, j) element of \mathbf{A} is denoted a_{ij} . An $(m \times n)$ all-ones

matrix where all of its elements are equal to 1 is denoted as $\mathbf{J}_{m \times n}$. An $(m \times n)$ zero matrix where all of its elements are equal to 0 is denoted as $\mathbf{O}_{m \times n}$. The $(n \times n)$ identity matrix is denoted \mathbf{I}_n . The k -dimensional column vector which has “1” in the j th row and zero elsewhere is called the unit vector and has denoted \mathbf{e}_j . The i th column of a matrix \mathbf{A} can be written as:

$$\mathbf{A}_{*i} = \mathbf{A} \cdot \underset{(q)}{\mathbf{e}_i}. \quad (44)$$

The elementary matrix

$$\mathbf{E}_{ij}^{(p \times q)} \triangleq \underset{(p)}{\mathbf{e}_i} \underset{(q)}{\mathbf{e}_j}^T \quad (45)$$

has dimensions $(p \times q)$, has “1” in the (i, j) th element, and has zero elsewhere.

The Kronecker product of a matrix \mathbf{A} with dimensions $(p \times q)$ and a matrix \mathbf{B} with dimensions $(m \times n)$ is represented as $\mathbf{A} \otimes \mathbf{B}$. The resulting matrix is of size $pm \times qn$ and is defined by

$$\mathbf{A} \otimes \mathbf{B} \triangleq \begin{bmatrix} a_{11}\mathbf{B} & a_{12}\mathbf{B} & \cdots & a_{1q}\mathbf{B} \\ a_{21}\mathbf{B} & & & \vdots \\ \vdots & & & \\ a_{p1}\mathbf{B} & \cdots & & a_{pq}\mathbf{B} \end{bmatrix}. \quad (46)$$

The Hadamard product of \mathbf{A} $(p \times q)$ and \mathbf{C} $(p \times q)$ is denoted $\mathbf{A} \odot \mathbf{C}$ and is a $(p \times q)$ matrix defined by

$$\mathbf{A} \odot \mathbf{C} \triangleq \begin{bmatrix} a_{11}c_{11} & a_{12}c_{12} & \cdots & a_{1q}c_{1q} \\ a_{21}c_{21} & & & \vdots \\ \vdots & & & \\ a_{p1}c_{p1} & \cdots & & a_{pq}c_{pq} \end{bmatrix}. \quad (47)$$

The permutation matrix is a square $(pq \times pq)$ matrix defined by [25]

$$\mathbf{U}_{p \times q} \triangleq \sum_i^p \sum_j^q \mathbf{E}_{ij}^{(p \times q)} \otimes \mathbf{E}_{ji}^{(q \times p)}. \quad (48)$$

Each row and column of the square matrix $\mathbf{U}_{p \times q}$ has only one element with a value of 1, and all remaining elements are zero [26].

A related matrix is a rectangular $(p^2 \times q^2)$ matrix defined by [25]

$$\bar{\mathbf{U}}_{p \times q} \triangleq \sum_i^p \sum_j^q \mathbf{E}_{ij}^{(p \times q)} \otimes \mathbf{E}_{ij}^{(p \times q)}. \quad (49)$$

The matrix derivative of a matrix function $\mathbf{F} = [f_{ij}(\mathbf{W})]_{m \times n} : \mathbb{R}^{p \times q} \rightarrow \mathbb{R}^{m \times n}$ with respect to a matrix $\mathbf{X}(p \times q)$ is of order $(pm \times qn)$ and is defined as [27]

$$\frac{\partial \mathbf{F}(\mathbf{X})}{\partial \mathbf{X}} \triangleq \begin{bmatrix} \frac{\partial \mathbf{F}(\mathbf{X})}{\partial x_{11}} & \cdots & \frac{\partial \mathbf{F}(\mathbf{X})}{\partial x_{1q}} \\ \vdots & & \vdots \\ \frac{\partial \mathbf{F}(\mathbf{X})}{\partial x_{p1}} & \cdots & \frac{\partial \mathbf{F}(\mathbf{X})}{\partial x_{pq}} \end{bmatrix}, \quad (50)$$

where

$$\frac{\partial \mathbf{F}(\mathbf{X})}{\partial x} \triangleq \begin{bmatrix} \frac{\partial f_{11}(\mathbf{X})}{\partial x} & \cdots & \frac{\partial f_{1n}(\mathbf{X})}{\partial x} \\ \vdots & & \vdots \\ \frac{\partial f_{m1}(\mathbf{X})}{\partial x} & \cdots & \frac{\partial f_{mn}(\mathbf{X})}{\partial x} \end{bmatrix} \quad (51)$$

for $x \in \mathbb{R}$.

A.2 Properties of Kronecker product and matrix calculus

The matrix operations related to Kronecker product are listed below and are adopted from [25]. The matrices in this subsection have the following dimension: $\mathbf{A} : \mathbb{R}^{s \times t} \rightarrow \mathbb{R}^{p \times q}$, $\mathbf{B} \in \mathbb{R}^{s \times t}$, and $\mathbf{C} : \mathbb{R}^{s \times t} \rightarrow \mathbb{R}^{q \times r}$.

$$\frac{\partial \mathbf{A}(\mathbf{B})\mathbf{C}(\mathbf{B})}{\partial \mathbf{B}} = \frac{\partial \mathbf{A}(\mathbf{B})}{\partial \mathbf{B}}(\mathbf{I}_t \otimes \mathbf{C}(\mathbf{B})) + (\mathbf{I}_s \otimes \mathbf{A}(\mathbf{B}))\frac{\partial \mathbf{C}(\mathbf{B})}{\partial \mathbf{B}}. \quad (52)$$

$$\frac{\partial \mathbf{A}}{\partial \mathbf{A}} = \bar{\mathbf{U}}_{p \times q}. \quad (53)$$

$$\frac{\partial \mathbf{A}^T}{\partial \mathbf{A}} = \mathbf{U}_{p \times q}. \quad (54)$$

The equations (52), (53), and (54) are employed in deriving the derivative of the loss function with respect to the weight matrix of GCN.

B Back-propagation of Graph Convolutional Network

A basic graph structure is defined as:

$$G = (V, E), \quad (55)$$

where $|V| = n$ is the number of nodes in the graph and $|E| = n_e$ is the number of edges. Denoting $v_i \in V$ as a node and $e_{ij} = (v_i, v_j) \in E$ as an edge pointing from v_i to v_j , the adjacency matrix, $\mathbf{A} \in \{0, 1\}^{n \times n}$, is an $n \times n$ matrix where a_{ij} equals 1 if the edge e_{ij} exists, and a_{ij} equals 0 if e_{ij} does not belong to E ; in addition, a graph may possess node attributes represented by the matrix $\mathbf{H}_0 \in \mathbb{R}^{n \times n_0}$, where $\mathbf{h}_{0_v} \in \mathbb{R}^{n_0}$ is the feature vector of a node v with n_0 features [20].

C Proof of theorem 2.3

Let $\mathbf{F}(\mathbf{W}) = [f_{ij}(\mathbf{W})]_{m \times n} : \mathbb{R}^{p \times q} \rightarrow \mathbb{R}^{m \times n}$ be a $m \times n$ multivariate matrix-valued function of a $p \times q$ matrix $\mathbf{W} \in \mathbb{R}^{p \times q}$, from definition 2.1, $\Sigma(\mathbf{F}(\mathbf{W}))$ maps a $m \times n$ matrix $\mathbf{F}(\mathbf{W})$ to a $m \times n$ matrix $\Sigma(\mathbf{F}(\mathbf{W}))$. From (50), the derivative of $\Sigma(\mathbf{F}(\mathbf{W}))$ ($m \times n$) with respect to \mathbf{W} ($p \times q$) can be written as:

$$\frac{\partial \Sigma(\mathbf{F}(\mathbf{W}))}{\partial \mathbf{W}} = \begin{bmatrix} \frac{\partial \Sigma(\mathbf{F}(\mathbf{W}))}{\partial w_{11}} & \dots & \frac{\partial \Sigma(\mathbf{F}(\mathbf{W}))}{\partial w_{1q}} \\ \vdots & & \vdots \\ \frac{\partial \Sigma(\mathbf{F}(\mathbf{W}))}{\partial w_{p1}} & \dots & \frac{\partial \Sigma(\mathbf{F}(\mathbf{W}))}{\partial w_{pq}} \end{bmatrix}. \quad (56)$$

Using (51), chain rule, and definition 2.2, the first element in (56) is given by

$$\begin{aligned} & \frac{\partial \Sigma(\mathbf{F}(\mathbf{W}))}{\partial w_{11}} \\ &= \begin{bmatrix} \frac{\partial \sigma(f_{11}(\mathbf{W}))}{\partial w_{11}} & \dots & \frac{\partial \sigma(f_{1n}(\mathbf{W}))}{\partial w_{11}} \\ \vdots & & \vdots \\ \frac{\partial \sigma(f_{m1}(\mathbf{W}))}{\partial w_{11}} & \dots & \frac{\partial \sigma(f_{mn}(\mathbf{W}))}{\partial w_{11}} \end{bmatrix} \\ &= \begin{bmatrix} \sigma'(f_{11}(\mathbf{W})) \frac{\partial f_{11}(\mathbf{W})}{\partial w_{11}} & \dots & \sigma'(f_{1n}(\mathbf{W})) \frac{\partial f_{1n}(\mathbf{W})}{\partial w_{11}} \\ \vdots & & \vdots \\ \sigma'(f_{m1}(\mathbf{W})) \frac{\partial f_{m1}(\mathbf{W})}{\partial w_{11}} & \dots & \sigma'(f_{mn}(\mathbf{W})) \frac{\partial f_{mn}(\mathbf{W})}{\partial w_{11}} \end{bmatrix} \\ &= \Sigma'(\mathbf{F}(\mathbf{W})) \odot \frac{\partial \mathbf{F}(\mathbf{W})}{\partial w_{11}}, \end{aligned} \quad (57)$$

where $\Sigma'(\mathbf{F}(\mathbf{W}))$ is defined as

$$\Sigma'(\mathbf{F}(\mathbf{W})) = \begin{bmatrix} \sigma'(f_{11}(\mathbf{W})) & \dots & \sigma'(f_{1n}(\mathbf{W})) \\ \vdots & & \vdots \\ \sigma'(f_{m1}(\mathbf{W})) & \dots & \sigma'(f_{mn}(\mathbf{W})) \end{bmatrix}. \quad (58)$$

Applying the calculation of (57) to all the other elements in (56), we can get

$$\begin{aligned}
& \frac{\partial \Sigma(\mathbf{F}(\mathbf{W}))}{\partial \mathbf{W}} \\
&= \begin{bmatrix} \Sigma'(\mathbf{F}(\mathbf{W})) \odot \frac{\partial \mathbf{F}(\mathbf{W})}{\partial w_{11}} & \cdots & \Sigma'(\mathbf{F}(\mathbf{W})) \odot \frac{\partial \mathbf{F}(\mathbf{W})}{\partial w_{1q}} \\ \vdots & & \vdots \\ \Sigma'(\mathbf{F}(\mathbf{W})) \odot \frac{\partial \mathbf{F}(\mathbf{W})}{\partial w_{p1}} & \cdots & \Sigma'(\mathbf{F}(\mathbf{W})) \odot \frac{\partial \mathbf{F}(\mathbf{W})}{\partial w_{pq}} \end{bmatrix} \\
&= \begin{bmatrix} \Sigma'(\mathbf{F}(\mathbf{W})) & \cdots & \Sigma'(\mathbf{F}(\mathbf{W})) \\ \vdots & & \vdots \\ \Sigma'(\mathbf{F}(\mathbf{W})) & \cdots & \Sigma'(\mathbf{F}(\mathbf{W})) \end{bmatrix} \odot \begin{bmatrix} \frac{\partial \mathbf{F}(\mathbf{W})}{\partial w_{11}} & \cdots & \frac{\partial \mathbf{F}(\mathbf{W})}{\partial w_{1q}} \\ \vdots & & \vdots \\ \frac{\partial \mathbf{F}(\mathbf{W})}{\partial w_{p1}} & \cdots & \frac{\partial \mathbf{F}(\mathbf{W})}{\partial w_{pq}} \end{bmatrix} \\
&= (\mathbf{J}_{p \times q} \otimes \Sigma'(\mathbf{F}(\mathbf{W}))) \odot \frac{\partial \mathbf{F}(\mathbf{W})}{\partial \mathbf{W}}.
\end{aligned} \tag{59}$$

D Back-propagation for multi-layer GCN with arbitrary activation functions

We consider a d -layer GCN defined as

$$\hat{\mathbf{Y}} = \Sigma_{d+1}(\mathbf{H}_d \mathbf{H}_d^T), \tag{60}$$

where

$$\mathbf{H}_d = \Sigma_d(\hat{\mathbf{A}} \mathbf{H}_{d-1} \mathbf{W}_d), \tag{61}$$

$d \in \mathbb{Z}^+$, $\hat{\mathbf{A}} = \tilde{\mathbf{D}}^{-\frac{1}{2}}(\mathbf{A} + \mathbf{I}_n)\tilde{\mathbf{D}}^{-\frac{1}{2}} \in \mathbb{R}^{n \times n}$ represents the $n \times n$ normalized adjacency matrix [28], $\tilde{\mathbf{D}} \in \mathbb{R}^{n \times n}$ is a degree matrix defined by $\tilde{d}_{ii} = 1 + \sum_j a_{ij}$, $\mathbf{A} \in \{0, 1\}^{n \times n}$ is an $n \times n$ adjacency matrix whose diagonal elements are all equal to zero, $\mathbf{H}_d \in \mathbb{R}^{n \times n_d}$ is the feature matrix for n nodes with n_d features, $\hat{\mathbf{Y}} \in \mathbb{R}^{n \times n_d}$ is the output matrix, and $\mathbf{W}_1 \in \mathbb{R}^{n_0 \times n_1}$, $\mathbf{W}_2 \in \mathbb{R}^{n_1 \times n_2}$, \dots , and $\mathbf{W}_d \in \mathbb{R}^{n_{d-1} \times n_d}$ are trainable parameter matrices, and $\Sigma_1, \Sigma_2, \dots$, and Σ_{d+1} are any element-wise activation function.

The loss function of the GCN for link prediction can be defined as

$$L = - \sum_{(i,j) \in E} \ln(\hat{y}_{ij}) - \sum_{(i,j) \in S} \ln(1 - \hat{y}_{ij}), \tag{62}$$

where y_{ij} denotes the element in the i th row and j th column of the training matrix $\mathbf{Y} \in \mathbb{R}^{n \times n_d}$, \hat{y}_{ij} denotes the element in the i th row and j th column of $\hat{\mathbf{Y}}$, and S is a set of edges that contains n_e negative edges $(\bar{v}_i, \bar{v}_j) \in S$ randomly sampled from E^c .

The derivative of the loss in (62) with respect to the s th weight matrix \mathbf{W}_s is

$$\frac{\partial L}{\partial \mathbf{W}_s} = - \sum_{(i,j) \in E} \frac{1}{\hat{y}_{ij}} \frac{\partial \hat{y}_{ij}}{\partial \mathbf{W}_s} + \sum_{(i,j) \in S} \frac{1}{1 - \hat{y}_{ij}} \frac{\partial \hat{y}_{ij}}{\partial \mathbf{W}_s}, \tag{63}$$

where $s \in \mathbb{Z}^+$, $s \leq d$,

$$\begin{aligned} & \frac{\partial \hat{y}_{ij}}{\partial \mathbf{W}_s} \\ &= \Sigma'_{d+1}(\mathbf{H}_{d_{i*}} \mathbf{H}_{d_{j*}}^T) \frac{\partial \mathbf{H}_{d_{i*}} \mathbf{H}_{d_{j*}}^T}{\partial \mathbf{W}_s} \end{aligned} \quad (64)$$

$$= \Sigma'_{d+1}(\mathbf{H}_{d_{i*}} \mathbf{H}_{d_{j*}}^T) \cdot \frac{\partial \Sigma_d(\hat{\mathbf{A}}_{i*} \mathbf{H}_{d-1} \mathbf{W}_d) \Sigma_d(\hat{\mathbf{A}}_{j*} \mathbf{H}_{d-1} \mathbf{W}_d)^T}{\partial \mathbf{W}_s} \quad (65)$$

$$\begin{aligned} &= \Sigma'_{d+1}(\mathbf{H}_{d_{i*}} \mathbf{H}_{d_{j*}}^T) \cdot \left(\frac{\partial \Sigma_d(\hat{\mathbf{A}}_{i*} \mathbf{H}_{d-1} \mathbf{W}_d)}{\partial \mathbf{W}_s} (\mathbf{I}_{n_s} \otimes \Sigma_d(\hat{\mathbf{A}}_{j*} \mathbf{H}_{d-1} \mathbf{W}_d)^T) \right. \\ &\quad \left. + (\mathbf{I}_{n_{s-1}} \otimes \Sigma_d(\hat{\mathbf{A}}_{i*} \mathbf{H}_{d-1} \mathbf{W}_d)) \frac{\partial \Sigma_d(\hat{\mathbf{A}}_{j*} \mathbf{H}_{d-1} \mathbf{W}_d)^T}{\partial \mathbf{W}_s} \right) \end{aligned} \quad (66)$$

$$\begin{aligned} &= \Sigma'_{d+1}(\mathbf{H}_{d_{i*}} \mathbf{H}_{d_{j*}}^T) \\ &\quad \cdot (((\mathbf{J}_{n_{s-1} \times n_s} \otimes \Sigma'_d(\hat{\mathbf{A}}_{i*} \mathbf{H}_{d-1} \mathbf{W}_d)) \odot \frac{\partial \hat{\mathbf{A}}_{i*} \mathbf{H}_{d-1} \mathbf{W}_d}{\partial \mathbf{W}_s}) \cdot (\mathbf{I}_{n_s} \otimes \Sigma_d(\hat{\mathbf{A}}_{j*} \mathbf{H}_{d-1} \mathbf{W}_d)^T) \\ &\quad + (\mathbf{I}_{n_{s-1}} \otimes \Sigma_d(\hat{\mathbf{A}}_{i*} \mathbf{H}_{d-1} \mathbf{W}_d)) \cdot ((\mathbf{J}_{n_{s-1} \times n_s} \otimes \Sigma'_d(\hat{\mathbf{A}}_{j*} \mathbf{H}_{d-1} \mathbf{W}_d)^T) \odot \frac{\partial (\hat{\mathbf{A}}_{j*} \mathbf{H}_{d-1} \mathbf{W}_d)^T}{\partial \mathbf{W}_s})), \end{aligned} \quad (67)$$

(64) follows from chain rule and definition 2.2, (65) follows from (61), (66) follows from (52), (67) follows from theorem (2.3),

$$\frac{\partial \hat{\mathbf{A}}_{i*} \mathbf{H}_{d-1} \mathbf{W}_d}{\partial \mathbf{W}_s} = \begin{cases} (\mathbf{I}_{n_{s-1}} \otimes (\hat{\mathbf{A}}_{i*} \mathbf{H}_{d-1})) \bar{\mathbf{U}}_{n_{s-1} \times n_s}, & \text{if } s = d \\ (\mathbf{I}_{n_{s-1}} \otimes \hat{\mathbf{A}}_{i*}) \frac{\partial \mathbf{H}_{d-1}}{\partial \mathbf{W}_s} (\mathbf{I}_{n_s} \otimes \mathbf{W}_d), & \text{if } s < d, \end{cases} \quad (68)$$

$$\frac{\partial (\hat{\mathbf{A}}_{j*} \mathbf{H}_{d-1} \mathbf{W}_d)^T}{\partial \mathbf{W}_s} = \begin{cases} \mathbf{U}_{n_{s-1} \times n_s} (\mathbf{I}_{n_s} \otimes (\hat{\mathbf{A}}_{j*} \mathbf{H}_{d-1})^T), & \text{if } s = d \\ (\mathbf{I}_{n_{s-1}} \otimes \mathbf{W}_d^T) \frac{\partial \mathbf{H}_{d-1}^T}{\partial \mathbf{W}_s} (\mathbf{I}_{n_s} \otimes \hat{\mathbf{A}}_{j*}^T), & \text{if } s < d, \end{cases} \quad (69)$$

(68) follows from (52) and (53), (69) follows from (52) and (54), $\frac{\partial \mathbf{H}_{d-1}}{\partial \mathbf{W}_s}$ is similar to (42) but with

$$\begin{aligned} & \mathbf{A} = \hat{\mathbf{A}}, \\ & \frac{\partial \mathbf{H}_{d-1}^T}{\partial \mathbf{W}_s} \\ &= \begin{cases} (\mathbf{J}_{n_{s-1} \times n_s} \otimes \Sigma'_{d-1}(\hat{\mathbf{A}} \mathbf{H}_{d-2} \mathbf{W}_{d-1})^T) \odot (\mathbf{U}_{n_{s-1} \times n_s} (\mathbf{I}_{n_s} \otimes (\hat{\mathbf{A}} \mathbf{H}_{d-2})^T)), & \text{if } s = d-1 \\ (\mathbf{J}_{n_{s-1} \times n_s} \otimes \Sigma'_{d-1}(\hat{\mathbf{A}} \mathbf{H}_{d-2} \mathbf{W}_{d-1})^T) \odot ((\mathbf{I}_{n_{s-1}} \otimes \mathbf{W}_{d-1}^T) \frac{\partial \mathbf{H}_{d-2}^T}{\partial \mathbf{W}_s} \cdot (\mathbf{I}_{n_s} \otimes \hat{\mathbf{A}}^T)), & \text{if } s < d-1, \end{cases} \end{aligned} \quad (70)$$

and (70) follows from (52), (54), and theorem (2.3).

E Sensitivity analysis

E.1 Sensitivity of loss with respect to feature matrix

The expression of the sensitivity of loss with respect to the feature matrix \mathbf{H}_0 is similar to the result in section 2.1.3 except that we change the dimension of the identity matrix, the all-one matrix, and the permutation related matrix, and we change the result in (35) and (42) for the last layer. Thus, the derivative, or sensitivity, of the loss in (39) with respect to the input feature matrix $\mathbf{H}_0 \in \mathbb{R}^{n \times n_0}$ is

$$\begin{aligned} \frac{\partial L}{\partial \mathbf{H}_0} &= - \sum_{i=1}^n \sum_{j=1}^{n_d} \left(\frac{y_{ij} - \hat{y}_{ij}}{\hat{y}_{ij}(1 - \hat{y}_{ij})} \Sigma'_{d+1}(h_{d_{ij}}) \cdot ((\mathbf{J}_{n \times n_0} \otimes \Sigma'_d(\mathbf{A}_{i*} \mathbf{H}_{d-1} \mathbf{W}_{d_{*j}})) \right. \\ &\quad \left. \odot \frac{\partial \mathbf{A}_{i*} \mathbf{H}_{d-1} \mathbf{W}_{d_{*j}}}{\partial \mathbf{H}_0}) \right), \end{aligned} \quad (71)$$

where $\frac{\partial \mathbf{A}_{i*} \mathbf{H}_{d-1} \mathbf{W}_{d*j}}{\partial \mathbf{H}_0}$ is defined as

$$\frac{\partial \mathbf{A}_{i*} \mathbf{H}_{d-1} \mathbf{W}_{d*j}}{\partial \mathbf{H}_0} = \begin{cases} (\mathbf{I}_n \otimes \mathbf{A}_{i*}) \cdot \bar{\mathbf{U}}_{n \times n_0} (\mathbf{I}_{n_0} \otimes \mathbf{W}_{d*j}) & \text{if } d = 1 \\ (\mathbf{I}_n \otimes \mathbf{A}_{i*}) \cdot \frac{\partial \mathbf{H}_{d-1}}{\partial \mathbf{H}_0} (\mathbf{I}_{n_0} \otimes \mathbf{W}_{d*j}) & \text{if } d > 1 \end{cases}, \quad (72)$$

and

$$\begin{aligned} & \frac{\partial \mathbf{H}_{d-1}}{\partial \mathbf{H}_0} \\ &= \begin{cases} (\mathbf{J}_{n \times n_0} \otimes \Sigma'_{d-1} (\mathbf{A} \mathbf{H}_{d-2} \mathbf{W}_{d-1})) \odot ((\mathbf{I}_n \otimes \mathbf{A}) \bar{\mathbf{U}}_{n \times n_0} (\mathbf{I}_{n_0} \otimes \mathbf{W}_{d-1})) & \text{if } d-1 = 1 \\ (\mathbf{J}_{n \times n_0} \otimes \Sigma'_{d-1} (\mathbf{A} \mathbf{H}_{d-2} \mathbf{W}_{d-1})) \odot ((\mathbf{I}_n \otimes \mathbf{A}) \frac{\partial \mathbf{H}_{d-2}}{\partial \mathbf{H}_0} (\mathbf{I}_{n_0} \otimes \mathbf{W}_{d-1})) & \text{if } d-1 > 1. \end{cases} \end{aligned} \quad (73)$$

E.2 Sensitivity of output with respect to feature matrix

Similarly, the derivative, or sensitivity, of each element of the output in (60) with respect to the input feature matrix $\mathbf{H}_0 \in \mathbb{R}^{n \times n_0}$ is

$$\begin{aligned} & \frac{\partial \hat{y}_{ij}}{\partial \mathbf{H}_0} \\ &= \Sigma'_{d+1} (\mathbf{H}_{d i*} \mathbf{H}_{d j*}^T) \\ & \quad \cdot (((\mathbf{J}_{n \times n_0} \otimes \Sigma'_d (\hat{\mathbf{A}}_{i*} \mathbf{H}_{d-1} \mathbf{W}_d)) \odot \frac{\partial \hat{\mathbf{A}}_{i*} \mathbf{H}_{d-1} \mathbf{W}_d}{\partial \mathbf{H}_0}) \cdot (\mathbf{I}_{n_0} \otimes \Sigma_d (\hat{\mathbf{A}}_{j*} \mathbf{H}_{d-1} \mathbf{W}_d)^T) \\ & \quad + (\mathbf{I}_n \otimes \Sigma_d (\hat{\mathbf{A}}_{i*} \mathbf{H}_{d-1} \mathbf{W}_d)) \cdot ((\mathbf{J}_{n \times n_0} \otimes \Sigma'_d (\hat{\mathbf{A}}_{j*} \mathbf{H}_{d-1} \mathbf{W}_d)^T) \odot \frac{\partial (\hat{\mathbf{A}}_{j*} \mathbf{H}_{d-1} \mathbf{W}_d)^T}{\partial \mathbf{H}_0})), \end{aligned} \quad (74)$$

where

$$\frac{\partial \hat{\mathbf{A}}_{i*} \mathbf{H}_{d-1} \mathbf{W}_d}{\partial \mathbf{H}_0} = \begin{cases} (\mathbf{I}_n \otimes \hat{\mathbf{A}}_{i*}) \bar{\mathbf{U}}_{n \times n_0} (\mathbf{I}_{n_0} \otimes \mathbf{W}_d), & \text{if } d = 1 \\ (\mathbf{I}_n \otimes \hat{\mathbf{A}}_{i*}) \frac{\partial \mathbf{H}_{d-1}}{\partial \mathbf{H}_0} (\mathbf{I}_{n_0} \otimes \mathbf{W}_d), & \text{if } d > 1, \end{cases} \quad (75)$$

$$\frac{\partial (\hat{\mathbf{A}}_{j*} \mathbf{H}_{d-1} \mathbf{W}_d)^T}{\partial \mathbf{H}_0} = \begin{cases} (\mathbf{I}_n \otimes \mathbf{W}_d^T) \mathbf{U}_{n \times n_0} (\mathbf{I}_{n_0} \otimes (\hat{\mathbf{A}}_{j*} \mathbf{H}_{d-1})^T), & \text{if } d = 1 \\ (\mathbf{I}_n \otimes \mathbf{W}_d^T) \frac{\partial \mathbf{H}_{d-1}^T}{\partial \mathbf{H}_0} (\mathbf{I}_{n_0} \otimes \hat{\mathbf{A}}_{j*}^T), & \text{if } d > 1, \end{cases} \quad (76)$$

$\frac{\partial \mathbf{H}_{d-1}}{\partial \mathbf{H}_0}$ is similar to (73) but with $\mathbf{A} = \hat{\mathbf{A}}$, and

$$\begin{aligned} & \frac{\partial \mathbf{H}_{d-1}^T}{\partial \mathbf{H}_0} \\ &= \begin{cases} (\mathbf{J}_{n \times n_0} \otimes \Sigma'_{d-1} (\hat{\mathbf{A}} \mathbf{H}_{d-2} \mathbf{W}_{d-1})^T) \odot ((\mathbf{I}_n \otimes \mathbf{W}_{d-1}^T) \mathbf{U}_{n \times n_0} (\mathbf{I}_{n_0} \otimes \hat{\mathbf{A}}^T)), & \text{if } d-1 = 1 \\ (\mathbf{J}_{n \times n_0} \otimes \Sigma'_{d-1} (\hat{\mathbf{A}} \mathbf{H}_{d-2} \mathbf{W}_{d-1})^T) \odot ((\mathbf{I}_n \otimes \mathbf{W}_{d-1}^T) \frac{\partial \mathbf{H}_{d-2}^T}{\partial \mathbf{H}_0} \cdot (\mathbf{I}_{n_0} \otimes \hat{\mathbf{A}}^T)), & \text{if } d-1 > 1. \end{cases} \end{aligned} \quad (77)$$

F Data description

F.1 Node classification

We tested our method for node classification on Zachary's Karate Club [19]. Zachary's Karate Club is a social network consisting of 34 members of a karate club, where undirected edges represent friendships, as shown in Figure 3 (a) [29]. Every node was labeled by one of four classes obtained via modularity-based clustering in [2] and we set class 2 to class 0 and class 3 to class 1 for binary node classification. The goal is to classify the nodes into the correct class.

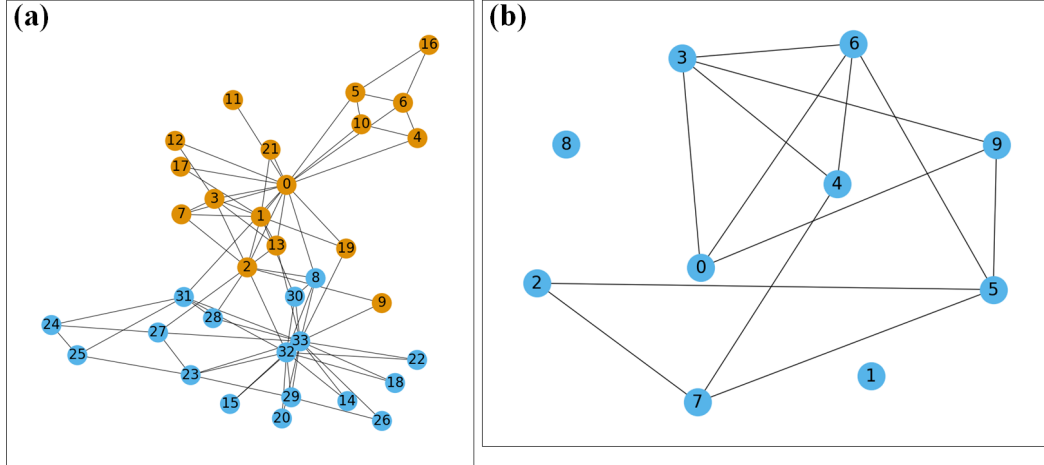


Figure 3: (a) Zachary’s karate club social network [19]. The node colors signify classes, with blue representing class 0 and orange representing class 1. (b) Drug-drug interaction network. Black links between each node represent graph edges.

F.2 Link prediction

Drug-Drug Interaction (DDI) commonly arises when multiple drugs that target the same receptors are administered concurrently. Such interactions can result in diminished therapeutic efficacy by hindering drug absorption and amplifying adverse drug reactions, which may lead to unforeseen off-target side effects [30]. We tested our method for link prediction using a DDI network, which was formed by extracting data from the DrugBank database [31]. This network revolved around drugs of interest, identified through the analysis of differentially expressed genes from RNA-sequence data within the GSE147507 dataset, particularly in the context of a COVID-19 study [32]. Interactions between drugs and proteins were established using DrugBank’s available drug-protein interactions, serving as communication pathways. The most disrupted signaling pathway, "Herpes simplex virus 1 infection," comprised 495 proteins and was selected based on the analysis of differentially expressed genes and KEGG pathway database [33]. Subsequently, 214 drugs were extracted from DrugBank for their direct interactions with the 495 proteins of interest. In total, 468 drugs (including the aforementioned 214) with 48,460 DDIs were sourced from DrugBank, capable of interacting with the 214 drugs. The DDI network was constructed using these 468 drugs. For this study, the first 100 drugs, based on their DrugBank index, were extracted to form a manageable subset for analysis.

The drug features analyzed in this study encompassed 20 molecular properties such as weight and water solubility. Furthermore, pairwise chemical similarities between drugs were computed using the Jaccard coefficient, quantifying structural similarity based on Simplified Molecular Input Line Entry System (SMILES) strings [34]. Principal Components Analysis was used and the first 20 (out of 488) principal components, which explained 98% of variability in drug feature data, were used as input features to GCN.

G The evolution of binary node classification and link prediction

G.0.1 5-layer GCN with identity function and sigmoid activation function for binary node classification

The 5-layer GCN is in the form of (37) and the loss function is defined by (39), where $d = 5$ is the number of layers, $\mathbf{H}_0 = \mathbf{I}_n \in \mathbb{R}^{n \times n}$ is the feature matrix for $n = 34$ nodes with $n = 34$ features, the matrix $\mathbf{Y} \in \mathbb{R}^{n \times n_d}$ represents the ground truth of nodes, $\hat{\mathbf{Y}} \in \mathbb{R}^{n \times n_d}$ denotes the final layer predictions of the GCN model, $n_d = 1$ for binary node classification, $\mathbf{W}_1 \in \mathbb{R}^{n_0 \times n_1}$, $\mathbf{W}_2 \in \mathbb{R}^{n_1 \times n_2}$, $\mathbf{W}_3 \in \mathbb{R}^{n_2 \times n_3}$, $\mathbf{W}_4 \in \mathbb{R}^{n_3 \times n_4}$, and $\mathbf{W}_5 \in \mathbb{R}^{n_4 \times n_d}$ are trainable parameter matrices, $n_0 = 20$, $n_1 = 2$, $n_2 = 3$, $n_3 = 2$, $n_4 = 3$, Σ_1 is an element-wise ReLU function [21], Σ_2 is an element-wise sigmoid linear Unit (SiLU) function [35], Σ_3 is an element-wise exponential linear unit (ELU)

function with $\alpha = 1$ [36], Σ_4 is an element-wise Leaky ReLU function [37], Σ_5 is an element-wise identity function, and Σ_6 is an element-wise sigmoid function.

The GCN is trained using SGD with a learning rate of 3×10^{-5} and with 10 iterations for both methods. The GCN is trained 1060 times and the weight matrices are reinitialized every time. If the calculation of the loss function has nan value, the training at this time will be skipped and will not be counted. The line plot of the SSE of the five weight matrices, \mathbf{W}_1 , \mathbf{W}_2 , \mathbf{W}_3 , \mathbf{W}_4 , and \mathbf{W}_5 , between the reverse mode automatic differentiation and our matrix-based method is shown in Figure 5 (a). The box plot of the \log_{10} of the SSE in Figure 5 (b) shows that the median of the five weight matrices is between 10^{-18} and 10^{-14} . In the box plot of Figure 5 (b), the box extends from the first quartile to the third quartile of the data, with a line at the median, while the horizontal lines extending from the box, which are called whiskers, indicating the farthest data point lying within 1.5 times the inter-quartile range from the box, and circles are those points past the end of the whiskers [38].

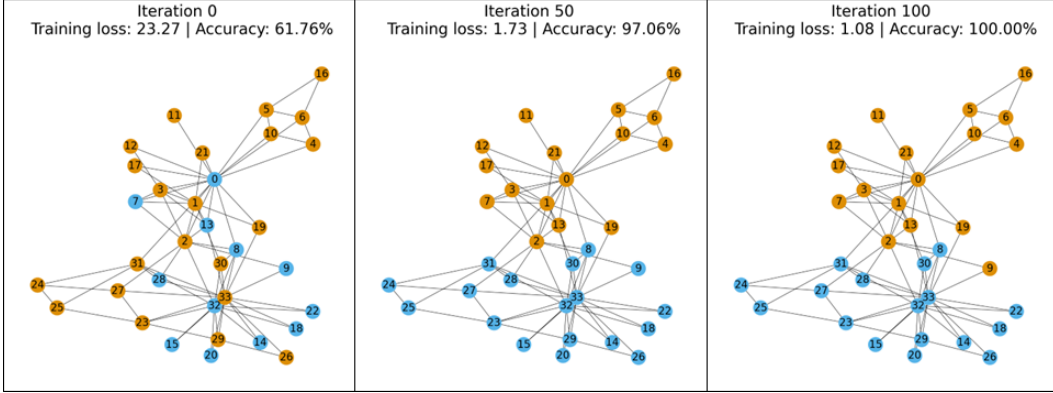


Figure 4: Evolution of karate club network node classification obtained from a 1-layer GCN model after 100 training iterations. The training loss, accuracy, and classification results exhibit similar trends when using either reverse-mode automatic differentiation or our matrix-based method. Nodes are represented by circle with the number to distinguish each node. Colors represent classes, where blue represents class 0 and orange represents class 1. Black links between each node represent graph edges.

G.0.2 5-layer GCN for link prediction

The 5-layer GCN is defined by (60) and the loss function is defined by (62), where $d = 5$ is the number of layers, the matrix $\hat{\mathbf{Y}} \in \mathbb{R}^{n \times n}$ represents ground-truth adjacency matrix, $\hat{\mathbf{Y}} \in \mathbb{R}^{n \times n}$ denotes the final link predictions of the GCN model, $\mathbf{W}_1 \in \mathbb{R}^{n_0 \times n_1}$, $\mathbf{W}_2 \in \mathbb{R}^{n_1 \times n_2}$, $\mathbf{W}_3 \in \mathbb{R}^{n_2 \times n_3}$, $\mathbf{W}_4 \in \mathbb{R}^{n_3 \times n_4}$, and $\mathbf{W}_5 \in \mathbb{R}^{n_4 \times n_5}$ are trainable parameter matrices, $n_0 = 20$, $n_1 = 2$, $n_2 = 3$, $n_3 = 5$, $n_4 = 3$, $n_5 = 40$, Σ_1 is an element-wise Leaky ReLU function [21], Σ_2 is an element-wise ELU function [35], Σ_3 is an element-wise SiLU function with $\alpha = 1$ [36], Σ_4 is an element-wise ReLU function [37], Σ_5 is an element-wise identity function, and Σ_6 is an element-wise sigmoid function.

The GCN is trained using SGD with a learning rate of 0.9 and with 10 iterations for both methods. In each iteration, thirteen negative edges are uniformly sampled from a set of all the unconnected edges for the GCN training using the reverse mode automatic differentiation. The sampled negative edges in each iteration are saved in a list and used in the training of the GCN using our matrix-based method. The GCN is trained 1060 times and the weight matrices are reinitialized every time. If the calculation of the loss function has nan value, the training at this time will be skipped and will not be counted. The line plot of the SSE of the five weight matrices, \mathbf{W}_1 , \mathbf{W}_2 , \mathbf{W}_3 , \mathbf{W}_4 , and \mathbf{W}_5 , between the reverse mode automatic differentiation and our matrix-based method is shown in Figure 6 (a). The box plot of the \log_{10} of the SSE in Figure 6 (b) shows that the median of the five weight matrices is between 10^{-18} and 10^{-14} .

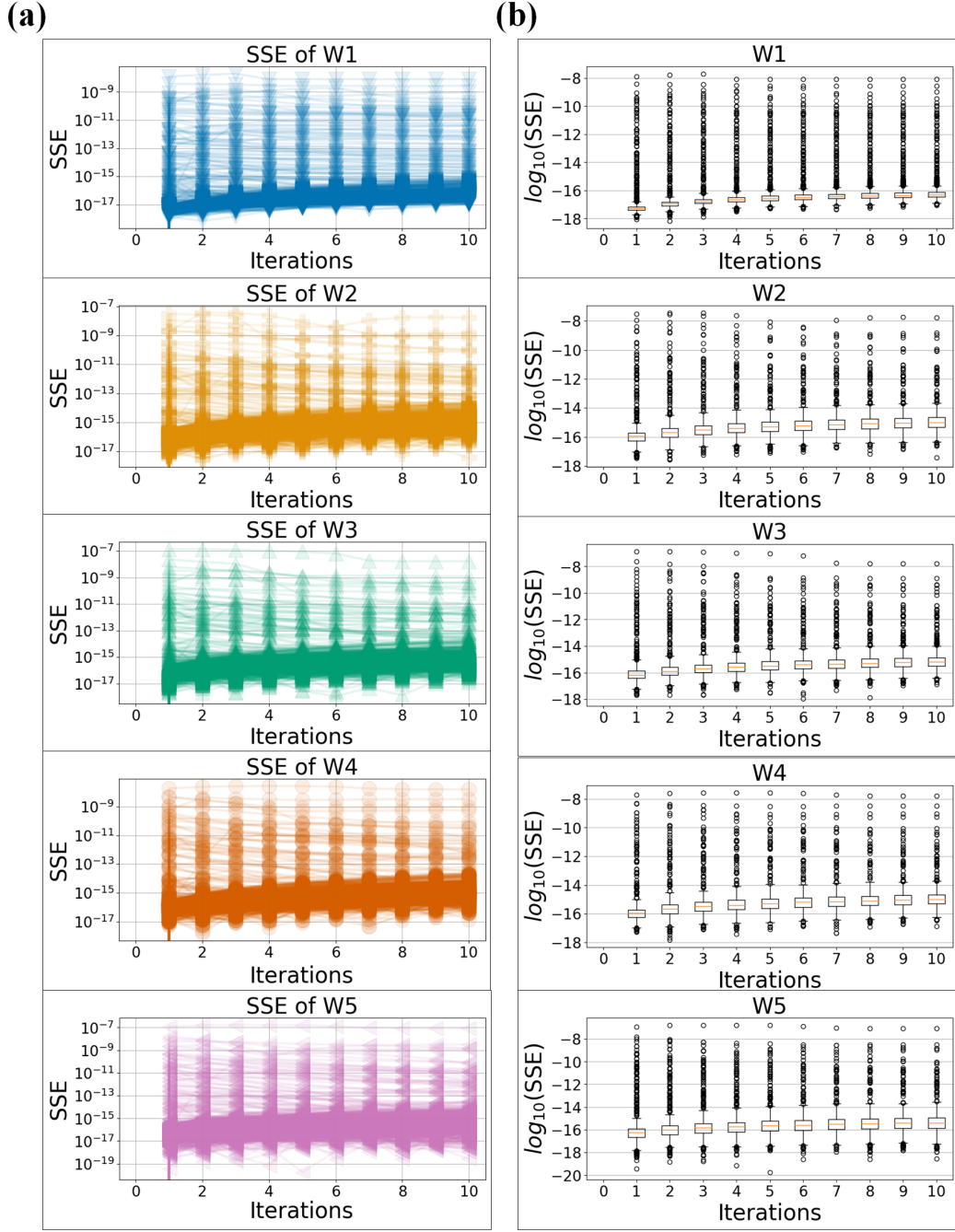


Figure 5: (a) The line plot of the 1060 evolution of the sum of squared error (SSE) between the trainable weight matrices obtained from our method and the matrix obtained using reverse mode automatic differentiation in section G.0.1. (b) The box plot of the 1060 evolution of the SSE between the trainable weight matrices obtained from our method and the matrix obtained using reverse mode automatic differentiation in section G.0.1.

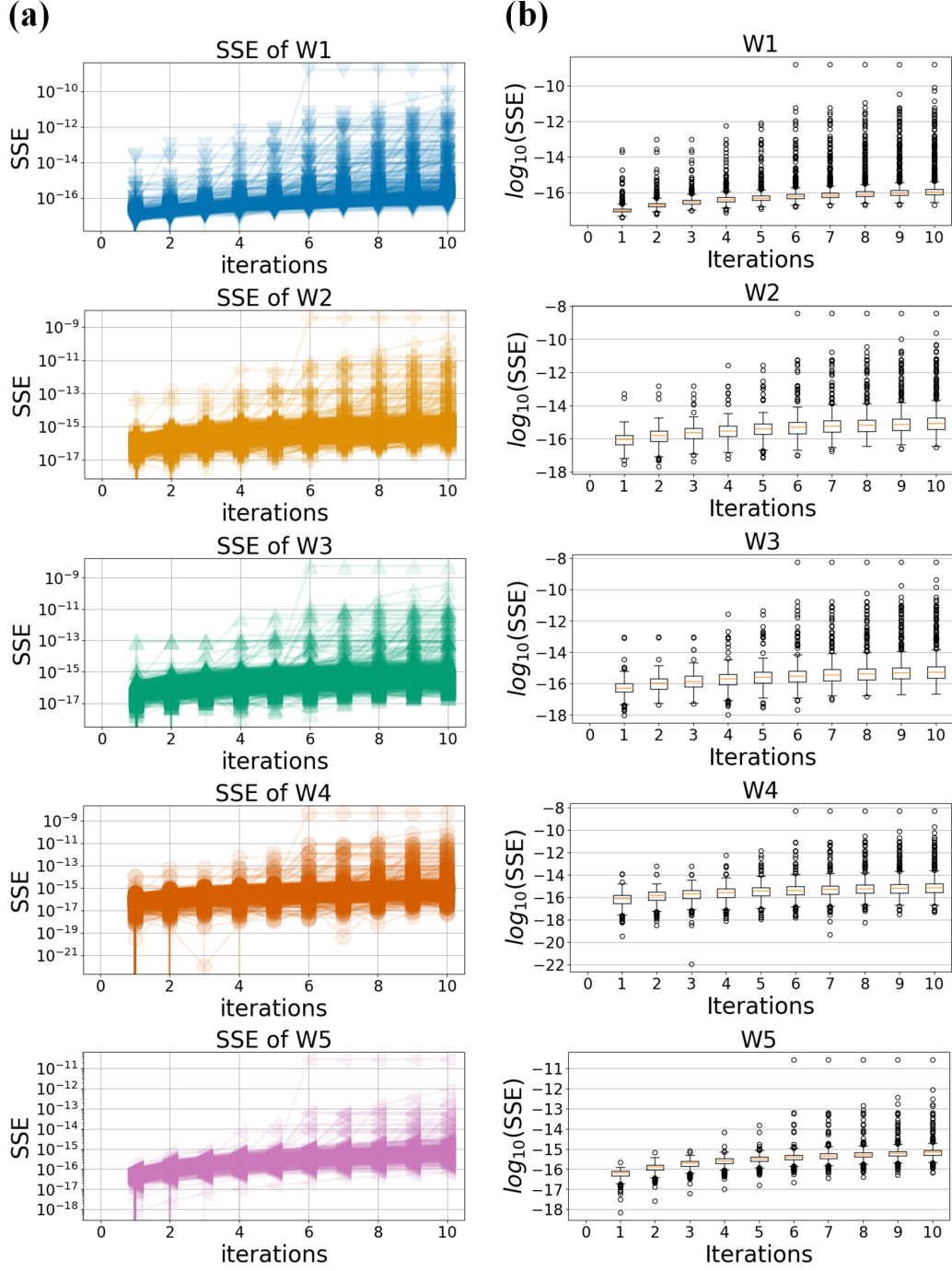


Figure 6: (a) The line plot of the 1060 evolution of the sum of squared error (SSE) between the trainable weight matrices obtained from our method and the matrix obtained using reverse mode automatic differentiation in section G.0.2. (b) The box plot of the 1060 evolution of the SSE between the trainable weight matrices obtained from our method and the matrix obtained using reverse mode automatic differentiation in section G.0.2.

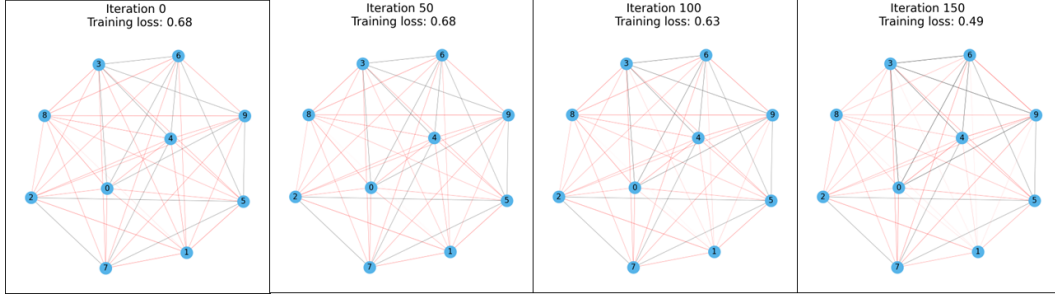


Figure 7: Evolution of DDI network node classification obtained from a 2-layer GCN model after 150 training iterations. The training loss and classification results exhibit similar trends when using either reverse-mode automatic differentiation or our matrix-based method. Nodes are represented by circle with the number to distinguish each node. Black links between each node represent truth graph edges. Red links between each node represent the edges that does not exist in the ground truth adjacency matrix. The transparency of the black and red link is proportional to the predicted link probability.



Figure 8: Heat maps of output sensitivity of the DDI network (1)

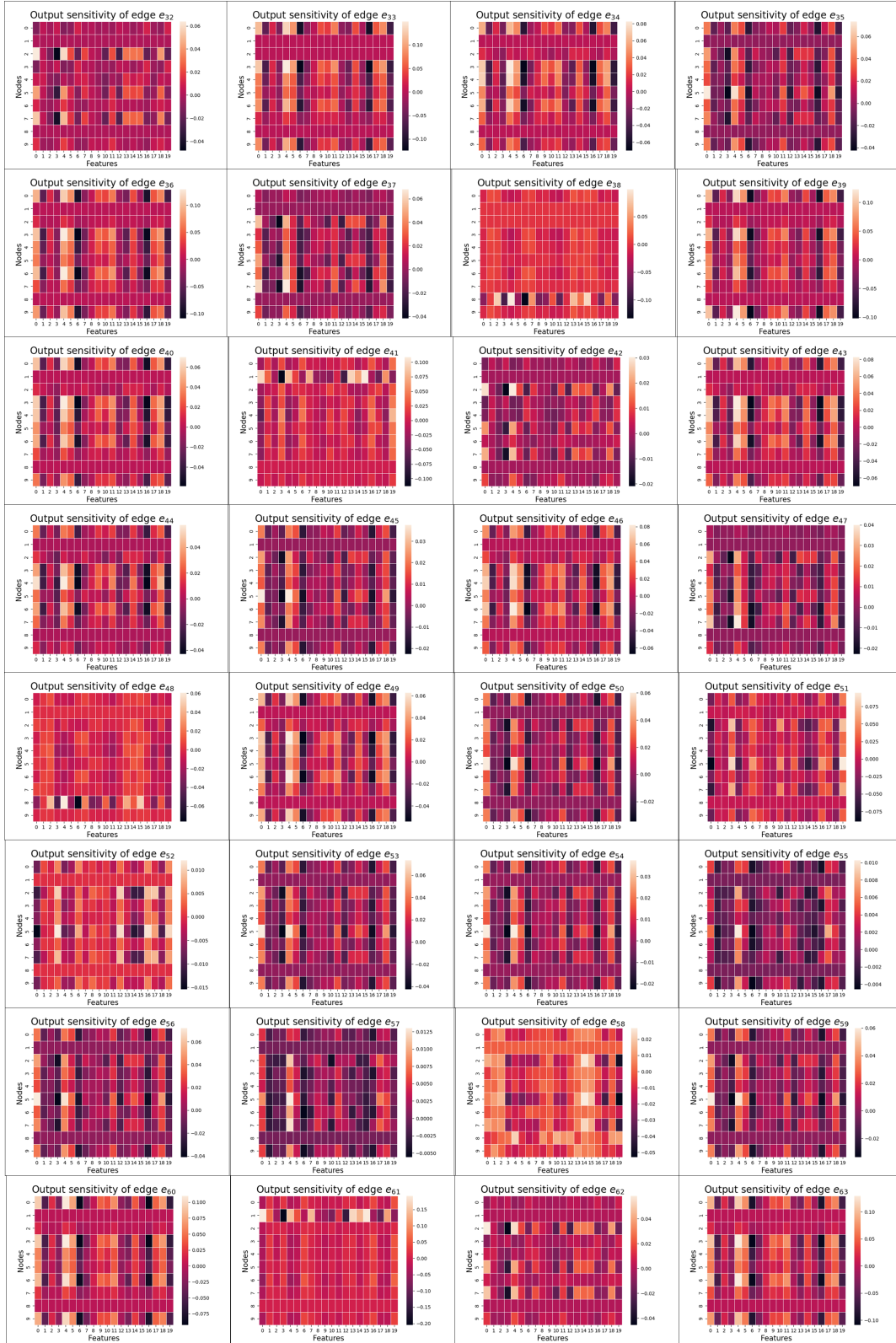


Figure 9: Heat maps of output sensitivity of the DDI network (2)

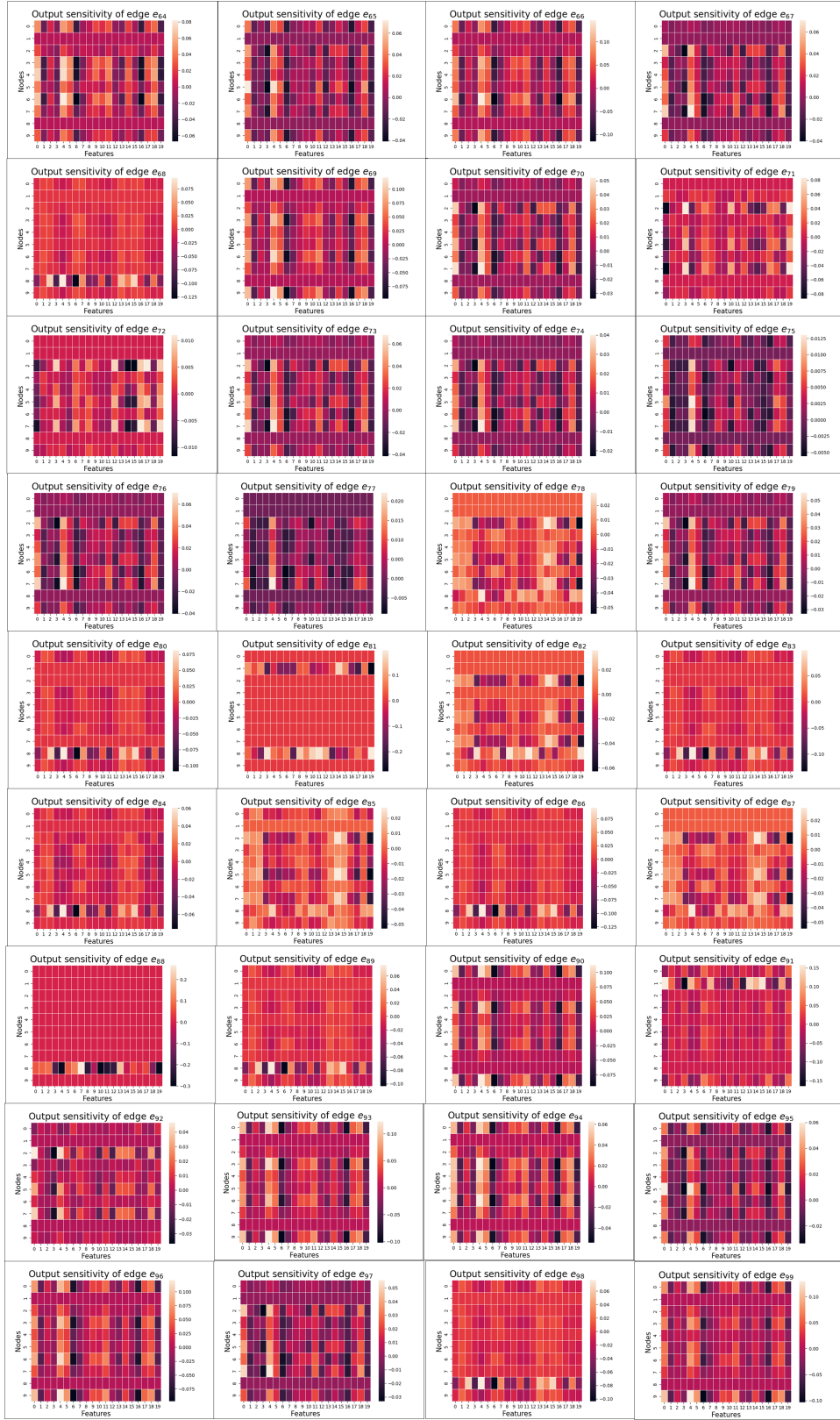


Figure 10: Heat maps of output sensitivity of the DDI network (3)