

---

# Towards Understanding the Generalization of Graph Neural Networks

---

Huayi Tang<sup>1,2</sup> Yong Liu<sup>1,2</sup>

## Abstract

Graph neural networks (GNNs) are the most widely adopted model in graph representation learning. Despite their extraordinary success in real-world applications, understanding their working mechanism by theory is still on primary stage. In this paper, we move towards this goal from the perspective of generalization. To be specific, we establish high probability bounds of generalization gap and gradients under transductive setting with consideration of stochastic optimization. After that, we provide high probability bounds of generalization gap for popular GNNs and discuss factors affecting their generalization capability. The theoretical results reveal the architecture specific factors affecting the generalization gap. Experimental results on benchmark datasets show the consistency between theoretical results and empirical evidence. Our results provide new insights in understanding the generalization of GNNs.

## 1. Introduction

Graph-structured data (Zhu et al., 2021) exists widely in real-world applications. As one of the most powerful tools to process graph-structured data, GNNs (Gori et al., 2005; Scarselli et al., 2009) are widely adopted in Computer Vision (Qi et al., 2017; Johnson et al., 2018; Landrieu & Simonovsky, 2018; Satorras & Estrach, 2018), Natural Language Processing (Bastings et al., 2017; Beck et al., 2018; Song et al., 2018), Recommendation Systems (Ying et al., 2018; Fan et al., 2019; He et al., 2020; Deng et al., 2022), AI for Science (Sanchez-Gonzalez et al., 2020; Pfaff et al., 2021; Shen et al., 2021; Han et al., 2022), to name a few. There are two main ways to view modern GNNs, *i.e.*, spatial domain perspective (Kipf & Welling, 2017; Veličković

et al., 2018; Xu et al., 2018; 2019) and spectral domain perspective (Defferrard et al., 2016; Gastegger et al., 2019; Liao et al., 2019; Chien et al., 2021; He et al., 2021). The former regards GNN as the process of combining and updating features according to adjacent relationships. The latter treats GNN as a filtering function applied on input features. Recent developments of GNNs are summarized in (Zhou et al., 2020; Wu et al., 2021; Zhang et al., 2022).

Despite the empirical success of GNNs, establishing theories to explain their behaviors is still in its infancy. Recent works towards this direction includes understanding over-smoothing (Li et al., 2018; Zhao & Akoglu, 2020; Oono & Suzuki, 2020a; Rong et al., 2020), interpretability (Ying et al., 2019; Luo et al., 2020; Vu & Thai, 2020; Yuan et al., 2020; 2021), expressiveness (Xu et al., 2019; Chen et al., 2019; Maron et al., 2019; Dehmamy et al., 2019; Feng et al., 2022), and generalization (Scarselli et al., 2018; Du et al., 2019; Verma & Zhang, 2019; Garg et al., 2020; Zhang et al., 2020; Oono & Suzuki, 2020b; Lv, 2021; Liao et al., 2021; Esser et al., 2021; Cong et al., 2021). This work focuses on the last branch. Some previous works adopt the classical techniques such as Vapnik-Chervonenkis dimension (Scarselli et al., 2018), Rademacher complexity (Lv, 2021; Garg et al., 2020) and algorithm stability (Verma & Zhang, 2019) to provide generalization bounds for GCN (Kipf & Welling, 2017) and more general message passing neural networks. However, in their analysis, the original graph is split into subgraphs composed of central node and its neighbors, which are treated as independent samples. This setting significantly differs from real implementation that training nodes are sampled without replacement from full nodes and the test nodes are visible during training (El-Yaniv & Pechyony, 2007; Oono & Suzuki, 2020a), resulting a gap between theory and practice. To tackle this issue, recent works (Oono & Suzuki, 2020b; Esser et al., 2021) incorporate the learning schema of GNNs into the category of transductive learning and derive more realistic results. However, there are still some drawbacks of these works. First, the analysis in (Oono & Suzuki, 2020b) is oriented to multi-scale GNNs that differ a lot from modern GNNs in network architecture. Besides, their analysis is limited to the AdaBoost-like optimization procedure, and whether the technique can be applied to general optimization algorithms such as stochastic gradient descent (SGD) is unknown. Second, the upper

---

<sup>1</sup>Gaoling School of Artificial Intelligence, Renmin University of China, Beijing, China <sup>2</sup>Beijing Key Laboratory of Big Data Management and Analysis Methods, Beijing, China. Correspondence to: Yong Liu <liuyonggsai@ruc.edu.cn>.

bound in (Esser et al., 2021) is of slow order and fails to provide meaningful learning guarantee for node classification in large-scale scenarios. Third, (Cong et al., 2021) only consider spectral-based GNNs with fixed coefficients, leaving spectral-based GNNs with learnable coefficients (Chien et al., 2021) unexplored.

Motivated by the aforementioned challenges, under transductive setting, we study the generalization gap of GNNs for node classification task with consideration of stochastic optimization algorithm. First, we establish high probability bounds of generalization gap and gradients under transductive setting, and derive high probability bounds of test error under gradient dominant condition. Next, we provide a comprehensive analysis on popular GNNs including both linear and non-linear models and derive the upper bound of the Lipschitz continuity and Hölder smoothness constants, by which we compare their generalization capability. The results show that SGC (Wu et al., 2019) and APPNP (Gasteiger et al., 2019) can achieve smaller generalization gap than GCN (Kipf & Welling, 2017). Besides, the unconstrained coefficients in spectral GNNs may lead to a large generalization gap. Our results reveal why shallow models yield comparable and even superior performance from the perspective of learning theory, and provide theoretical supports for widely used techniques such as early stop and drop edge (Rong et al., 2020). Experimental results on benchmark datasets show that the theoretical findings are generally consistent with the practical evidences.

## 2. Related Work

### 2.1. Generalization Analysis of GNNs

Existing studies on the generalization of GNNs general fall into two categories: graph classification task and node classification task.

**Graph classification task.** (Liao et al., 2021) is the first work to establish generalization bounds of GCN and message passing neural networks by PAC-Bayesian approach. The authors in (Ju et al., 2023) further improve their results and provide the lower bound. Besides, neural tangent kernels (Jacot et al., 2018) are also used to analyze the generalization of infinitely wide GNNs trained by gradient descent (Du et al., 2019). Different from that, this work focus on node classification task that is more challenging.

**Node classification task.** The authors in (Scarselli et al., 2018) analyze the generalization capability of GNNs by Vapnik–Chervonenkis dimension. (Verma & Zhang, 2019) is the first work to provide generalization bounds of one-layer GCN by algorithm stability which is further extended to multi-layer GCNs in (Zhou & Wang, 2021). The work (Garg et al., 2020) converts the graph into individual local node-wise computation tree and bound their generalization

bound respectively by Rademacher Complexity. The aforementioned works rely on the assumption that converting a graph into subgraphs, which differs a lot from realistic implementation. Observing that, (Oono & Suzuki, 2020b) makes the first step that adopting the transductive learning framework to analyze multi-scale GNNs. This framework originates from (Vapnik, 1998; 2006), and is further developed in (El-Yaniv & Pechyony, 2006; 2007) where the authors propose transductive stability and Transductive Rademacher complexity to measure the generalization capability of transductive learner. The work most related to ours is (Cong et al., 2021) and (Esser et al., 2021), where the authors establish generalization bound for GNNs and its variants by transductive uniform stability and transductive Rademacher complexity respectively. However, the derived bound in (Esser et al., 2021) is of slow order, and whether their technique can be applied on SGD is still unknown. Different from (Cong et al., 2021) that analyzing full-batch gradient descent, we analyze a more complex setting, *i.e.*, transductive learning under SGD, due to the involve of randomness in optimization. Besides, there are some works orthogonal to ours, *e.g.*, analyzing the generalization capability of GNNs training with topology-sampling (Li et al., 2022a) or on large random graphs (Keriven et al., 2020).

### 2.2. Out-of-Distribution (OOD) Generalization on Graphs

Much efforts are devoted to the study of OOD generalization on graphs (Li et al., 2022b) in recent years, due to the occurs of distribution shift in real-world scenarios. An adversarial learning schema (Wu et al., 2022) is proposed to minimize the mean and variance of risks from multiple environments. The authors in (Yang et al., 2022) propose a two-stage training schema to tackle distribution shift on molecular graphs. Energy-based message passing scheme is show to be effective in enhancing the OOD detection performance of GNNs (Wu et al., 2023). Current work (Yang et al., 2023) shows that the spurious performance of GNNs may come from its intrinsic generalization capability rather than expressivity. Besides, there are also some work focus on the reasoning (Xu et al., 2020), extrapolation ability (Xu et al., 2021; Bevilacqua et al., 2021), and generalization from small to large graphs (Yehudai et al., 2021).

## 3. Preliminaries

### 3.1. Notations

Let  $\mathcal{G} = \{\mathcal{V}, \mathcal{E}\}$  be an given undirected graph with  $n = |\mathcal{V}|$  nodes. Each node is an instance  $z_i = (\mathbf{x}_i, y_i)$  containing feature  $\mathbf{x}_i$  and label  $y_i$  from some space  $\mathcal{Z} = \mathcal{X} \times \mathcal{Y}$ . Let  $\mathbf{X}$  be the feature matrix where the  $i$ -th row  $\mathbf{X}_{i*}$  is the node feature  $\mathbf{x}_i$ . Let  $\mathbf{A}$  and  $\mathbf{D}$  be the adjacency matrix and the diagonal degree matrix respectively, where  $\mathbf{D}_{ii} = \sum_{j=1}^n \mathbf{A}_{ij}$ .

Denote by  $\tilde{\mathbf{A}} = (\mathbf{D} + \mathbf{I}_n)^{-\frac{1}{2}}(\mathbf{A} + \mathbf{I}_n)(\mathbf{D} + \mathbf{I}_n)^{-\frac{1}{2}}$  the normalized adjacency matrix with self-loops and  $\sqrt{|\mathcal{Y}|}$  the number of categories. We focus on the transductive learning setting in this work, *i.e.*, all features together with the randomly sampled labels are constructed as training set. Let  $S = \{\mathbf{x}_i, y_i\}_{i=1}^{m+u}$  be the set of instances where  $m + u = n$ . Without loss of generality (w.l.o.g.), let  $\{y_i\}_{i=1}^m$  be the selected labels, our task is to predict the labels of samples  $\{\mathbf{x}_i\}_{i=m+1}^{m+u}$  by a learner (model) trained on  $\{\mathbf{x}_i\}_{i=1}^{m+u} \cup \{y_i\}_{i=1}^m$ . This setting is widely adopted in node classification task (Yang et al., 2016; Kipf & Welling, 2017) where the training and test nodes are determined by a random partition.

From now on, we limit the scope of the learner to a given GNN and let  $\{\mathbf{W}_h\}_{h=1}^H$  be its learnable parameters. Since  $\mathbb{R}^{p \times q}$  and  $\mathbb{R}^{pq}$  are isomorphic, the analysis in this work is oriented to the vector space for concise. To this end, we use a unified vector  $\mathbf{w} = [\text{vec}[\mathbf{W}_1]; \dots; \text{vec}[\mathbf{W}_H]]$  to represent the collection of  $\{\mathbf{W}_h\}_{h=1}^H$ , where  $\text{vec}[\cdot]$  is the vectorization operator that transforms a given matrix into vector, *i.e.*,  $\text{vec}[\mathbf{W}] = [\mathbf{W}_{*1}; \dots; \mathbf{W}_{*q}]$  for  $\mathbf{W} \in \mathbb{R}^{p \times q}$ . Here  $\mathbf{W}_{*i}$  is the  $i$ -th column of  $\mathbf{W}$ . For  $\mathbf{w} \in \mathcal{W}$ , the training and test error is defined as  $R_m(\mathbf{w}) \triangleq \frac{1}{m} \sum_{i=1}^m \ell(\mathbf{w}; z_i)$  and  $R_u(\mathbf{w}) \triangleq \frac{1}{u} \sum_{i=m+1}^{m+u} \ell(\mathbf{w}; z_i)$  respectively, where  $\ell : \mathcal{W} \times \mathcal{Z} \mapsto \mathbb{R}_+$  is the loss function. In this work, we follow previous studies (El-Yaniv & Pechyony, 2007; Oono & Suzuki, 2020b; Esser et al., 2021) and define the transductive generalization gap by  $|R_m(\mathbf{w}) - R_u(\mathbf{w})|$ . Since the label of test examples are not available, the optimization process is finding parameters to minimize the training error  $R_m(\mathbf{w})$ . Much efforts (Duchi et al., 2011; Kingma & Ba, 2015) are devoted to solve this stochastic optimization problem, and we mainly focus on SGD (Summarized in Algorithm 1) in this work.

Now we introduce notations used in the rest of this paper. Denote by  $\|\cdot\|_2$  and  $\|\cdot\|$  the 2-norm of vector and spectral norm of matrix, respectively. Let  $\mathbf{w}^{(1)}$  be the initialization weight of model, we focus on the space  $\mathcal{W} = B(\mathbf{w}^{(1)}; r)$ ,  $r \geq 1$  in this work, where  $B(\mathbf{w}^{(1)}; r) \triangleq \{\mathbf{w} : \|\mathbf{w} - \mathbf{w}^{(1)}\|_2 \leq r\}$  is the ball with radius  $r$ . Denote by  $\nabla \ell(\cdot; z)$  the gradient of  $\ell$  with respect to (w.r.t.) the first argument. Denote by  $b_g = \sup_{z \in \mathcal{Z}} \|\nabla \ell(\mathbf{w}^{(1)}; z)\|_2$  the supremum of gradient with initialed parameter and  $b_\ell = \sup_{z \in \mathcal{Z}} |\ell(\mathbf{w}^{(1)}; z)|_2$  the supremum of loss value with initialed parameter. Let  $\hat{\mathbf{w}} \in \text{argmin}_{\mathbf{w} \in \mathcal{W}} R_m(\mathbf{w})$  be the parameters of training error minimizer. We denote by  $\sigma(\cdot)$  the activation function.

### 3.2. Assumptions

In this part, we present the assumptions used in this paper.

**Assumption 3.1.** Assume that there exists a constant  $c_X > 0$  such that  $\|\mathbf{x}\|_2 \leq c_X$  holds for all  $\mathbf{x} \in \mathcal{X}$ .

---

#### Algorithm 1 SGD for Transductive Learning

---

**Input:** Initial parameter  $\mathbf{w}^{(1)}$ , learning rates  $\{\eta_t\}$ , training set  $\{\mathbf{x}_i\}_{i=1}^{m+u} \cup \{y_i\}_{i=1}^m$ .

**for**  $t = 1$  **to**  $T$  **do**

Randomly draw  $j_t$  from the uniform distribution over the set  $\{j : j \in [m]\}$ .

Update parameters by

$$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} - \eta_t \nabla \ell(\mathbf{w}^{(t)}; z_{j_t}).$$

**end for**

---

**Assumption 3.2.** Assume that there exists a constant  $c_W > 0$  such that  $\|\mathbf{W}_h\| \leq c_W$ ,  $h \in [H]$  for  $\mathbf{w} \in B(\mathbf{w}^{(1)}; r)$ .

*Remark 3.3.* Assumption 3.1 requires that input features are bounded (Verma & Zhang, 2019). This assumption can be satisfied by applying normalization on features. Assumption 3.2 means that the parameters during the training process are bounded, which is a common assumption in generalization analysis of GNNs (Garg et al., 2020; Liao et al., 2021; Cong et al., 2021; Esser et al., 2021). These two assumptions are necessary to analyze the Lipschitz continuity and Hölder smoothness of objective w.r.t.  $\mathbf{w}$ .

**Assumption 3.4.** Assume that the activation function  $\sigma(\cdot)$  is  $\tilde{\alpha}$ -Hölder smooth. To be specific, let  $P > 0$  and  $\tilde{\alpha} \in (0, 1]$ , for all  $\mathbf{u}, \mathbf{v} \in \mathbb{R}^d$ ,

$$\|\sigma'(\mathbf{u}) - \sigma'(\mathbf{v})\|_2 \leq P \|\mathbf{u} - \mathbf{v}\|_2^{\tilde{\alpha}}.$$

*Remark 3.5.* It can be verified that Assumption 3.4 implies Lipschitz continuity of activation function if  $\tilde{\alpha} = 0$ . Besides, Assumption 3.4 implies the smoothness of activation function if  $\tilde{\alpha} = 1$ . Therefore, Assumption 3.4 is much milder than the assumption in previous work (Verma & Zhang, 2019; Cong et al., 2021) that requires the activation function is smooth. For the convenience of analysis while not yielding a large gap between theory and practice, we construct a modified ReLU function (See Appendix A) with hyperparameter  $q \in (1, 2]$  that satisfies Assumption 3.4 and has a tolerable approximate error to vanilla ReLU function.

**Assumption 3.6.** Assume that there exist a constant  $G > 0$  such that for all  $z \in S$

$$\sqrt{\eta_t} \|\nabla \ell(\mathbf{w}_t; z)\|_2 \leq G$$

holds  $\forall t \in \mathbb{N}$ , where  $\{\eta_t\}_{t=1}^T$  is learning rates.

*Remark 3.7.* A formal definition of  $\nabla \ell(\mathbf{w}; z)$  is provided in Lemma A.4 in the Appendix. Assumption 3.6 (Lei & Tang, 2021; Li & Liu, 2021) means that the product of gradient and the square root of learning rate is bounded, which is milder than the widely used bounded gradient assumption (Hardt et al., 2016; Kuzborskij & Lampert, 2018), since the learning rate tends to zero during the iteration.

**Assumption 3.8.** Assume that there exists a constant  $\sigma_0 > 0$  such that for  $\forall t \in \mathbb{N}_+$ , the following inequality holds

$$\mathbb{E}_{j_t} [\|\nabla \ell(\mathbf{w}; z_{j_t})\|_2] \leq \sigma_0^2.$$

*Remark 3.9.* Assumption 3.8 requires the boundness of variances of stochastic gradients, which is a standard assumption in stochastic optimization studies (Kuzborskij & Lampert, 2018; Lei & Tang, 2021; Li & Liu, 2021).

## 4. Theoretical Results

In this section, we first present the high probability bounds of generalization gap and excess risks under transductive learning in Section 4.1. After that, we turn to specific examples and provide results of some popular GNNs in Section 4.2. Please refer to the Appendix for complete proofs.

### 4.1. General Results of Transductive SGD

We first analyze properties of the objective function  $\ell$  and provide the following proposition.

**Proposition 4.1** (Informal). *Suppose Assumptions 3.1, 3.2, and 3.4 hold. Denote by  $\mathcal{F}$  a specific GNN, for any  $\mathbf{w}, \mathbf{w}' \in \mathcal{W}$  and  $z \in S$ , the objective  $\ell(\mathbf{w}; z)$  satisfies*

$$|\ell(\mathbf{w}; z) - \ell(\mathbf{w}'; z)| \leq L_{\mathcal{F}} \|\mathbf{w} - \mathbf{w}'\|_2, \quad (1)$$

and

$$\begin{aligned} & \|\nabla \ell(\mathbf{w}; z) - \nabla \ell(\mathbf{w}'; z)\|_2 \\ & \leq P_{\mathcal{F}} \max \{ \|\mathbf{w} - \mathbf{w}'\|_2^{\tilde{\alpha}}, \|\mathbf{w} - \mathbf{w}'\|_2 \}, \end{aligned} \quad (2)$$

with constant  $L_{\mathcal{F}}$  and  $P_{\mathcal{F}}$ .

*Remark 4.2.* We provide more detailed analysis to  $L_{\mathcal{F}}$  and  $P_{\mathcal{F}}$  in Section 4.2. Both  $L_{\mathcal{F}}$  and  $P_{\mathcal{F}}$  depend on the specific network architecture  $\mathcal{F}$  of GNNs. Thus, the upper bound of generalization gap vary by the architecture.

Our first main result is high probability bounds on the transductive generalization gap, as presented in Theorem 4.3.

**Theorem 4.3.** *Suppose Assumptions 3.1, 3.2, 3.4, 3.6, and 3.8 hold. Suppose that the learning rate  $\{\eta_t\}$  satisfies  $\eta_t = \frac{1}{t+t_0}$  such that  $t_0 \geq \max\{(2P)^{1/\alpha}, 1\}$ . For any  $\delta \in (0, 1)$ , with probability  $1 - \delta$ ,*

(a). *If  $\alpha \in (0, \frac{1}{2})$ , we have*

$$\begin{aligned} & R_u(\mathbf{w}_1^{(T+1)}) - R_m(\mathbf{w}^{(T+1)}) \\ & = \mathcal{O}\left(L_{\mathcal{F}} \frac{(m+u)^{\frac{3}{2}}}{mu} \log^{\frac{1}{2}}(T) T^{\frac{1-2\alpha}{2}} \log\left(\frac{1}{\delta}\right)\right). \end{aligned}$$

(b). *If  $\alpha = \frac{1}{2}$ , we have*

$$\begin{aligned} & R_u(\mathbf{w}^{(T+1)}) - R_m(\mathbf{w}^{(T+1)}) \\ & = \mathcal{O}\left(L_{\mathcal{F}} \frac{(m+u)^{\frac{3}{2}}}{mu} \log(T) \log\left(\frac{1}{\delta}\right)\right). \end{aligned}$$

(c). *If  $\alpha \in (\frac{1}{2}, 1]$ , we have*

$$\begin{aligned} & R_u(\mathbf{w}^{(T+1)}) - R_m(\mathbf{w}^{(T+1)}) \\ & = \mathcal{O}\left(L_{\mathcal{F}} \frac{(m+u)^{\frac{3}{2}}}{mu} \log^{\frac{1}{2}}(T) \log\left(\frac{1}{\delta}\right)\right). \end{aligned}$$

*Remark 4.4.* Theorem 4.3 shows that the transductive generalization gap depends on the training/test data size  $m/u$ , network architecture related Lipschitz continuity constant  $L_{\mathcal{F}}$ , and the number of iterations  $T$ . Generally, our upper bounds are of order  $\mathcal{O}\left(\left(\frac{1}{m} + \frac{1}{u}\right)\sqrt{m+u}\right)$ , which is much sharper than the bound  $\mathcal{O}\left(\left(\frac{1}{m} + \frac{1}{u}\right)(m+u) + \log(m+u)\right)$  in previous work (Esser et al., 2021). Note that with the increase of data size  $m+u$ , the bound in (Esser et al., 2021) become increasing larger and fail to provide a reasonable generalization guarantee. This seriously restricts its application in large-scale node classification scenarios where the order of  $m+u$  is usually millions. Our results address these drawbacks and provide more applicable generalization guarantee for GNNs. Besides, the bound provided in (Esser et al., 2021) does not consider the specific optimization and has difficulty in revealing the influence of  $T$  on generalization gap. Our result shows that the generalization gap becomes larger when the number of  $T$  increases, resulting in the over-fitting phenomenon. Thus, early stop may be beneficial for yielding a smaller generalization gap, which is widely adopted in implementation of modern GNNs (Kipf & Welling, 2017; Chen et al., 2020). It can be seen that the generalization gap is positively related to the Lipschitz continuity constant  $L_{\mathcal{F}}$  determined by specific network architecture  $\mathcal{F}$ . Thus, larger  $L_{\mathcal{F}}$  leads to larger upper bounds of generalization gap, showing that the network architecture of GNN also have a significant influence on the generalization gap (See Section 4.2 for more detail). The upper bound of generalization gap in (Cong et al., 2021) also increase with  $T$  when the objective is optimized by full-batch gradient descent. This is not surprise since it can be seen as a special case of SGD where the batch size is equal to the size of training samples.

Our second main result is high probability bounds of the gradients on training and test data.

**Theorem 4.5.** *Suppose Assumptions 3.1, 3.2, 3.4, 3.6, and 3.8 hold. Suppose that the learning rate  $\{\eta_t\}$  satisfies  $\eta_t = \frac{1}{t+t_0}$  such that  $t_0 \geq \max\{(2P)^{1/\alpha}, 1\}$ . For any  $\delta \in (0, 1)$ , with probability  $1 - \delta$ ,*

(a). *If  $\alpha \in (0, \frac{1}{2})$ , we have*

$$\begin{aligned} & \left\| \nabla R_m(\mathbf{w}^{(T+1)}) - \nabla R_u(\mathbf{w}^{(T+1)}) \right\|_2 \\ & = \mathcal{O}\left(\frac{(m+u)^{\frac{3}{2}}}{mu} \log^{\frac{1}{2}}(T) T^{\frac{1-2\alpha}{2}} \log\left(\frac{1}{\delta}\right)\right). \end{aligned}$$

(b). If  $\alpha = \frac{1}{2}$ , we have

$$\begin{aligned} & \left\| \nabla R_m(\mathbf{w}^{(T+1)}) - \nabla R_u(\mathbf{w}^{(T+1)}) \right\|_2 \\ &= \mathcal{O} \left( \frac{(m+u)^{\frac{3}{2}}}{mu} \log(T) \log \left( \frac{1}{\delta} \right) \right). \end{aligned}$$

(c). If  $\alpha \in (\frac{1}{2}, 1]$ , we have

$$\begin{aligned} & \left\| \nabla R_m(\mathbf{w}^{(T+1)}) - \nabla R_u(\mathbf{w}^{(T+1)}) \right\|_2 \\ &= \mathcal{O} \left( \frac{(m+u)^{\frac{3}{2}}}{mu} \log^{\frac{1}{2}}(T) \log \left( \frac{1}{\delta} \right) \right). \end{aligned}$$

**Remark 4.6.** Theorem 4.5 provides high probability bounds for the generalization gap of gradients under transductive setting. Overall, the generalization gap we derive is still of order  $\mathcal{O} \left( \left( \frac{1}{m} + \frac{1}{u} \right) \sqrt{m+u} \right)$ , which is applicable in real-world large-scale graph dataset. Besides, the generalization gap of gradients increases with the increase of  $T$ , showing that a smaller number of iterations helps achieving a smaller generalization gap of gradients.

Since the generalization performance is determined by both training error and generalization gap, we provide an upper bound of the test error under a special case that the objective satisfies the following PL condition.

**Assumption 4.7.** Suppose that there exists a constant  $\mu$  such that for all  $\mathbf{w} \in \mathcal{W}$ ,

$$R_m(\mathbf{w}) - R_m(\hat{\mathbf{w}}^*) \leq \frac{1}{2\mu} \|\nabla R_m(\mathbf{w})\|_2,$$

holds for the given set  $S$  from  $\mathcal{Z}$ .

**Remark 4.8.** Assumption 4.7 is also named as gradient dominance condition in learning theory studies, indicating that the difference between the optimal training error and the current training error can be upper bounded by the quadratic function of the gradient on training instances. This assumption is widely adopted in nonconvex learning (Zhou et al., 2018; Xu & Zeevi, 2020; Lei & Tang, 2021; Li & Liu, 2021), and has been verified in over-parameterized systems including wide neural networks (Liu et al., 2020). This assumption only appears in Theorem 4.9.

**Corollary 4.9.** Suppose Assumptions 3.1, 3.2, 3.4, 3.6, 3.8, and 4.7 hold. Suppose that the learning rate  $\{\eta_t\}$  satisfies  $\eta_t = \frac{2}{\mu(t+t_0)}$  such that  $t_0 \geq \max\{\frac{2}{\mu}(2P)^{\frac{1}{\alpha}}, 1\}$ . For any  $\delta \in (0, 1)$ , with probability  $1 - \delta$ ,

(a). If  $\alpha \in (0, \frac{1}{2})$ , we have

$$\begin{aligned} & R_u(\mathbf{w}^{(T+1)}) - R_m(\mathbf{w}^*) \\ &= \mathcal{O} \left( L_{\mathcal{F}} \frac{(m+u)^{\frac{3}{2}}}{mu} \log^{\frac{1}{2}}(T) T^{\frac{1}{2}-\alpha} \log \left( \frac{1}{\delta} \right) + \frac{1}{T^\alpha} \right), \end{aligned}$$

(b). If  $\alpha = \frac{1}{2}$ , we have

$$\begin{aligned} & R_u(\mathbf{w}^{(T+1)}) - R_m(\mathbf{w}^*) \\ &= \mathcal{O} \left( L_{\mathcal{F}} \frac{(m+u)^{\frac{3}{2}}}{mu} \log(T) \log \left( \frac{1}{\delta} \right) + \frac{1}{T^\alpha} \right). \end{aligned}$$

(c). If  $\alpha \in (\frac{1}{2}, 1)$ , we have

$$\begin{aligned} & R_u(\mathbf{w}^{(T+1)}) - R_m(\mathbf{w}^*) \\ &= \mathcal{O} \left( L_{\mathcal{F}} \frac{(m+u)^{\frac{3}{2}}}{mu} \log^{\frac{1}{2}}(T) \log(1/\delta) + \frac{1}{T^\alpha} \right). \end{aligned}$$

(d). If  $\alpha = 1$ , we have

$$\begin{aligned} & R_u(\mathbf{w}^{(T+1)}) - R_u(\mathbf{w}^*) \\ &= \mathcal{O} \left( L_{\mathcal{F}} \frac{(m+u)^{\frac{3}{2}}}{mu} \log^{\frac{1}{2}}(T) \log(1/\delta) + \frac{\log(T) \log^3(1/\delta)}{T} \right). \end{aligned}$$

**Remark 4.10.** Theorem 4.9 shows that under Assumption 4.7, the test error are determined by the minimal training error, optimization error and generalization gap. The minimal training error reflects how well the model fits data, which is a measure of the expressive ability. The first and the second term in the slack terms are generalization gap and optimization error, respectively. With the increase of  $T$ , the generalization gap increase while the optimization error decrease. Therefore, it is necessary to carefully choose a proper number of iterations in order to balance the trade-off between optimization and generalization. In the implementation of most GNNs studies (Kipf & Welling, 2017; Veličković et al., 2018; Chien et al., 2021; He et al., 2021), early stop is widely adopted and  $T$  is determined by the performance of model on validation set. Thus, our results are consistent with real implementations.

It is worth point out that although the results in this section is oriented to the case that the objective has two parameters (e.g., GCN, APPNP, and GPR-GNN in Section 4.2), results for other cases that the objective has one parameter (e.g., SGC in Section 4.2) or three parameters (e.g., GCNII in Section 4.2) have the same form when neglecting the constant factors. Meanwhile, the assumptions need to be modified correspondingly. Readers are referred to the Appendix for detailed discussion.

## 4.2. Cases Study of Popular GNNs

We have established high probability bounds for transductive generalization gap in Theorem 4.3. In this part, we analyze the upper bounds of architecture related constant  $L_{\mathcal{F}}$  and  $P_{\mathcal{F}}$ , with that the upper bound of generalization gap can be determined. Five representative GNNs, including GCN, GCNII, SGC, APPNP, and GPR-GNN, are selected for analysis. The loss function  $\ell$  is cross-entropy loss and denote

by  $\hat{\mathbf{Y}}$  the prediction. For concise, we do not consider the bias term, since it can be verified that  $\langle \mathbf{w}, \mathbf{x} \rangle + b = \langle \tilde{\mathbf{w}}, \tilde{\mathbf{x}} \rangle$  holds with  $\tilde{\mathbf{w}} = [\mathbf{w}; b]$  and  $\tilde{\mathbf{x}} = [\mathbf{x}; 1]$ .

**GCN.** The work (Kipf & Welling, 2017) proposes to aggregate features from one-hop neighbor nodes. The feature propagation process of a two-layer GCN model is

$$\hat{\mathbf{Y}} = \text{Softmax}(g(\tilde{\mathbf{A}})\sigma(g(\tilde{\mathbf{A}})\mathbf{X}\mathbf{W}_1)\mathbf{W}_2), \quad (3)$$

where  $g(\tilde{\mathbf{A}}) = \tilde{\mathbf{A}}$  and  $\mathbf{W}_1 \in \mathbb{R}^{d \times h}$ ,  $\mathbf{W}_2 \in \mathbb{R}^{h \times |\mathcal{Y}|}$  are parameters.

**Proposition 4.11.** *Suppose Assumptions 3.1, 3.2, and 3.4 hold, then the objective  $\ell(\mathbf{w}; z)$  is  $L_{\mathcal{F}}$ -Lipschitz continuous and Hölder smooth w.r.t.  $\mathbf{w} = [\text{vec}[\mathbf{W}_1]; \text{vec}[\mathbf{W}_2]]$ . Concretely, the Lipschitz continuity constant  $L_{\mathcal{F}}$  is  $L_{\text{GCN}} = 2c_X c_W \|\tilde{\mathbf{A}}\|_{\infty}^2$ .*

Due to the tedious formulation, we provide the concrete value of  $P_{\mathcal{F}}$  in the Appendix. Proposition 4.11 demonstrates that  $L_{\text{GCN}}$  mainly depends on factors  $\|g(\tilde{\mathbf{A}})\|_{\infty}$ ,  $c_X$ , and  $c_W$ . Let  $\text{deg}_{\min}$  and  $\text{deg}_{\max}$  be the minimum and maximum node degree, respectively. By Lemma A.1 in Appendix A,

$$\|\tilde{\mathbf{A}}\|_{\infty} \leq \sqrt{\frac{\text{deg}_{\max} + 1}{\text{deg}_{\min} + 1}}. \quad (4)$$

It can be found that the generalization gap decreases with the decrease of the maximum node degree, which could be achieved by removing edges. This explains in some sense why the DropEdge (Rong et al., 2020) technique is beneficial for alleviating the over-fitting problem from the perspective of learning theory. Besides, for GCN trained on sampled sub-graphs  $\{\mathcal{G}_i\}_{i=1}^n$ , the Lipschitz continuity constant is  $L_{\text{GCN}} = 2c_X c_W \max_{i \in [n]} \|\tilde{\mathbf{A}}^{[i]}\|_{\infty}^2$ , where  $\tilde{\mathbf{A}}^{[i]}$  is the normalized adjacency matrix with self-loop of  $\mathcal{G}_i$ . Since only a portion of neighboring nodes are preserved during sub-graphs sampling (Hamilton et al., 2017; Zeng et al., 2020; 2021), the maximum node degree of each sub-graph is smaller than that of initial graph, implying  $\max_{i \in [n]} \|\tilde{\mathbf{A}}^{[i]}\|_{\infty} \leq \|\tilde{\mathbf{A}}\|_{\infty}$  holds. Thus, Proposition 4.11 shows that training on sampled sub-graphs are beneficial to achieve smaller generalization gap. Lastly, the spectral norm of learning parameters also has an effect on the generalization gap. Thus, the commonly used  $L_2$  regularization technique is beneficial to reduce the generalization gap.

**GCNII.** The authors in (Chen et al., 2020) propose to relieve over-smoothing by initial residual and identity mapping. Denote by  $\mathbf{H}^{(0)} = \sigma(\mathbf{X}\mathbf{W}_0)$  the initial representation. The forward propagation of a two-layer GCNII model is

$$\begin{aligned} \mathbf{H}^{(1)} &= \sigma\left(\left((1 - \alpha_1)g(\tilde{\mathbf{A}})\mathbf{H}^{(0)} + \alpha_1\mathbf{H}^{(0)}\right)\Psi(\beta_1, \mathbf{W}_1)\right), \\ \mathbf{H}^{(2)} &= \sigma\left(\left((1 - \alpha_2)g(\tilde{\mathbf{A}})\mathbf{H}^{(1)} + \alpha_2\mathbf{H}^{(0)}\right)\Psi(\beta_2, \mathbf{W}_2)\right), \\ \hat{\mathbf{Y}} &= \text{softmax}(\mathbf{H}^{(2)}\mathbf{W}_3), \end{aligned}$$

where  $\Psi(\beta, \mathbf{W}) = (1 - \beta)\mathbf{I} + \beta\mathbf{W}$  and  $g(\tilde{\mathbf{A}}) = \tilde{\mathbf{A}}$ .  $\mathbf{W}_1 \in \mathbb{R}^{d \times h}$ ,  $\mathbf{W}_2 \in \mathbb{R}^{h \times h}$ , and  $\mathbf{W}_3 \in \mathbb{R}^{h \times |\mathcal{Y}|}$  are parameters.

**Proposition 4.12.** *Suppose Assumptions 3.1, 3.2, and 3.4 hold, then the objective  $\ell(\mathbf{w}; z)$  is  $L_{\mathcal{F}}$  Lipschitz continuous and Hölder smooth w.r.t.*

$$\mathbf{w} = [\text{vec}[\mathbf{W}_0]; \text{vec}[\mathbf{W}_1]; \text{vec}[\mathbf{W}_2]; \text{vec}[\mathbf{W}_3]].$$

Specifically, denote by  $C_{\ell} = 1 - \beta_{\ell} + \beta_{\ell}c_W$ ,  $\ell \in [2]$  and

$$\begin{aligned} B_1 &= c_X c_W C_1 \left( (1 - \alpha_1) \|\tilde{\mathbf{A}}\|_{\infty} + \alpha_1 \right), \\ B_2 &= \left( (1 - \alpha_2) B_1 \|\tilde{\mathbf{A}}\|_{\infty} + \alpha_2 c_X c_W \right) C_2, \\ L_1 &= 2 \left( 2 + \frac{c_W^2 \beta_2^2}{C_2^2} \right) B_2^2, \\ L_2 &= 2(1 - \alpha_2)^2 \beta_1^2 c_W^2 \|\tilde{\mathbf{A}}\|_{\infty}^2 \left( \frac{B_1^2 C_2^2}{C_1^2} \right). \end{aligned} \quad (5)$$

The Lipschitz continuity constant is  $L_{\text{GCNII}} = \sqrt{L_1 + L_2}$ .

Proposition 4.12 shows that  $L_{\text{GCNII}}$  is a function of  $\{\alpha_i\}_{i=1}^2$  and  $\{\beta_i\}_{i=1}^2$ . Finding the optimal value of  $L_{\text{GCNII}}$  is a quadratic programming problem with constrain  $\alpha_1, \alpha_2 \in [0, 1]$  and  $\beta_1, \beta_2 \in [0, 1]$ . Now we discuss a special case that  $\alpha_1 = \alpha_2 = 0$  and  $\beta_1 = \beta_2 = 0$ . In this case, we have  $L_1 = 4c_X^2 c_W^2 \|\tilde{\mathbf{A}}\|_{\infty}^4$  and  $L_2 = 0$ , which implies that  $L_{\text{GCNII}} = L_{\text{GCN}}$ . Note that the optimal value of  $L_{\text{GCNII}}$  is no larger than any value of objective function over the feasible region. Therefore, we conclude that the value of  $L_{\text{GCNII}}$  is no higher than  $L_{\text{GCN}}$ . This result is not surprise, since GCNII is a special GCN model under this setting. For proper value of  $\{\alpha_i\}_{i=1}^2$  and  $\{\beta_i\}_{i=1}^2$ , GCNII could achieve smaller generalization gap than GCN. As GCNII can achieve lower training error by relieving the over-smoothing problem, Proposition 4.12 indicates that GCNII can achieve superior performance when hyperparameters are set properly. Due to the involve of  $\{\alpha_i\}_{i=1}^2$  and  $\{\beta_i\}_{i=1}^2$ , the growth rate of  $L_{\text{GCNII}}$  is much smaller than  $L_{\text{GCN}}$  when propagation depth increases, which makes GCNII maintain generalization capability and achieve stale performance (See Section 5).

**SGC.** The work (Wu et al., 2019) proposes to remove all the nonlinear activation in GCN. To facilitate comparison with GCN, we consider a two layers SGC model, whose propagation is given by

$$\hat{\mathbf{Y}} = \text{softmax}(g(\tilde{\mathbf{A}})\mathbf{X}\mathbf{W}_1\mathbf{W}_2), \quad (6)$$

where  $g(\tilde{\mathbf{A}}) = \tilde{\mathbf{A}}^2$ .  $\mathbf{W}_1 \in \mathbb{R}^{d \times h}$  and  $\mathbf{W}_2 \in \mathbb{R}^{h \times |\mathcal{Y}|}$  is the parameter.

**Proposition 4.13.** *Suppose Assumption 3.1, 3.2, and 3.4 hold, then the objective  $\ell(\mathbf{w}; z)$  is  $L_{\mathcal{F}}$ -Lipschitz continuous and Hölder smooth w.r.t.  $\mathbf{w} = [\text{vec}[\mathbf{W}_1]; \text{vec}[\mathbf{W}_2]]$ . Specifically, the Lipschitz continuity constant  $L_{\mathcal{F}}$  is  $L_{\text{SGC}} = 2c_X c_W \|\tilde{\mathbf{A}}^2\|_{\infty}$ .*

Since  $\|\tilde{\mathbf{A}}^2\|_\infty \leq \|\tilde{\mathbf{A}}\|_\infty^2$ , we have  $L_{\text{SGC}} \leq L_{\text{GCN}}$ . Surprisingly, this simple linear model can achieve better smaller generalization gap than those nonlinear models (Kipf & Welling, 2017; Chen et al., 2020; Gasteiger et al., 2019; Chien et al., 2021), even though its representation ability is inferior than them. Note that the performance on test samples is determined by both training error and generalization gap. If linear GNNs can achieve a small training error, it is natural that they can achieve comparable and even better performance than nonlinear GNNs on test samples. Therefore, Proposition 4.13 reveals why linear GNNs achieve better performance than nonlinear GNNs from learning theory, as observed in recent works (Wu et al., 2019; Zhu & Koniusz, 2021; Wang et al., 2021). Considering the efficiency and scalability of linear GNNs on large-scale datasets, we believe that they have much potential to be exploited.

**APPNP.** Multi-scale features are aggregated via personalized PageRank schema in (Gasteiger et al., 2019). Formally, the feature propagation process is formulated as

$$\hat{\mathbf{Y}} = \text{softmax}(g(\tilde{\mathbf{A}})\sigma(\sigma(\mathbf{X}\mathbf{W}_1)\mathbf{W}_2)), \quad (7)$$

where  $g(\tilde{\mathbf{A}}) = \sum_{k=0}^{K-1} \gamma(1-\gamma)^k \tilde{\mathbf{A}}^k + (1-\gamma)^K \tilde{\mathbf{A}}^K$ .  $\mathbf{W}_1 \in \mathbb{R}^{d \times h}$  and  $\mathbf{W}_2 \in \mathbb{R}^{h \times |\mathcal{Y}|}$  are the parameters.

**Proposition 4.14.** *Suppose Assumption 3.1, 3.2, and 3.4 hold, then the objective  $\ell(\mathbf{w}; z)$  is  $L_{\mathcal{F}}$ -Lipschitz continuous and Hölder smooth w.r.t.  $\mathbf{w} = [\text{vec}[\mathbf{W}_1]; \text{vec}[\mathbf{W}_2]]$ . Concretely, the Lipschitz continuity constant  $L_{\mathcal{F}}$  is  $L_{\text{APPNP}} = 2c_X c_W \|g(\tilde{\mathbf{A}})\|_\infty$ .*

The Lipschitz continuity constant in Proposition 4.14 is positively related to the infinity matrix norm of the polynomial spectral filter. According to (Gasteiger et al., 2019),  $\gamma$  is commonly set to be a small number, yielding that  $\|g(\tilde{\mathbf{A}})\|_\infty < \|\tilde{\mathbf{A}}\|_\infty$  holds. Thus, the Lipschitz continuity constant of APPNP is smaller than that of GCN, indicating that APPNP may achieve smaller generalization gap than GCN. Besides,  $K$  also affects the value of  $\|g(\tilde{\mathbf{A}})\|_\infty$ , and a larger  $K$  may yield a larger generalization gap. Therefore,  $K$  is usually set as a proper value to guarantee a trade-off between expressive ability and generalization performance.

**GPR-GNN.** Compared with APPNP, the fixed coefficients are replaced by learnable weights in (Chien et al., 2021), in order to adaptively simulate both high-pass and low-pass graph filters. The feature propagation process is

$$\hat{\mathbf{Y}} = (g(\tilde{\mathbf{A}}, \gamma)\sigma(\sigma(\mathbf{X}\mathbf{W}_1)\mathbf{W}_2)), \quad (8)$$

where  $g(\tilde{\mathbf{A}}, \gamma) = \sum_{k=0}^K \gamma_k \tilde{\mathbf{A}}^k$ .  $\mathbf{W}_1 \in \mathbb{R}^{d \times h}$ ,  $\mathbf{W}_2 \in \mathbb{R}^{h \times |\mathcal{Y}|}$  and  $\gamma \in \mathbb{R}^{K+1}$  are the parameters.

**Proposition 4.15.** *Suppose Assumption 3.1, 3.2, and 3.4 hold, then the objective  $\ell(\mathbf{w}; z)$  is  $L_{\mathcal{F}}$ -Lipschitz continuous and Hölder smooth w.r.t.  $\mathbf{w} = [\text{vec}[\mathbf{W}_1]; \text{vec}[\mathbf{W}_2]; \gamma]$ .*

Concretely, the Lipschitz continuity constant  $L_{\mathcal{F}}$  is  $L_{\text{GPR}} = \sqrt{L_1^2 + L_2^2}$ , where

$$\begin{aligned} L_1 &= \sqrt{2}c_X c_W^2 \left( \sum_{k=0}^K \|\tilde{\mathbf{A}}^k\|_\infty \right), \\ L_2 &= 2c_X c_W \|g(\tilde{\mathbf{A}}, \gamma)\|_\infty. \end{aligned} \quad (9)$$

Note that  $L_2$  has similar form with  $L_{\text{APPNP}}$  (the only difference lie on the definition of  $g(\tilde{\mathbf{A}}, \gamma)$ ). Assume that  $g(\tilde{\mathbf{A}}, \gamma) = g(\tilde{\mathbf{A}})$  and note that  $L_{\text{GPR}} = \sqrt{L_1^2 + L_2^2} \geq L_2$ , we have  $L_{\text{GPR}} \geq L_{\text{APPNP}}$ . Besides, since there is no constraint on  $\gamma$ , the value of  $\|g(\tilde{\mathbf{A}}, \gamma)\|_\infty$  may be larger when the norm of  $\gamma$  is large, resulting in larger generalization gap than APPNP. Therefore, adopting regularization technique on the learnable coefficients to restrict the value of  $\|g(\tilde{\mathbf{A}}, \gamma)\|_\infty$  is necessary.

To summarize,  $L_{\mathcal{F}}$  and  $P_{\mathcal{F}}$  are determined by the feature propagation process and graph-structured data. Estimating these constants precisely is challenging (Virmaux & Scaman, 2018; Fazlyab et al., 2019), and the upper bounds we provided are sufficient to reflect the realistic generalization gap of these models (See Section 5 for more detail). Besides, we have to emphasize that results for GCN and GCNII with more than two layers can be derived by similar techniques, yet it requires more tedious computation. Exploring new techniques to estimate these constants conveniently and precisely are left for future work.

## 5. Experiments

**Experimental Setup.** We conduct experiments on widely adopted benchmark datasets, including Cora, Citeseer, and Pubmed (Sen et al., 2008; Yang et al., 2016). The accuracy and loss gap (*i.e.*, the absolute value of difference between the loss (accuracy) on training and test samples) are used to estimate the generalization gap. Following the standard transductive learning setting, in each run, 30% sampled nodes determined by a random seed are used as training set and the rest nodes are treated as test set. The number of iterations is fixed to  $T = 300$ . We independently repeat the experiments for 10 times and report the mean value and standard deviations of all runs. Please see Appendix C for more detailed settings.

**Experimental Results.** The loss and accuracy comparisons are presented in Table 3 and Table 1, respectively. We have the following observations: (1) SGC and APPNP have smaller loss and accuracy gap than other model including GCN, which is consistent with the analysis in Proposition 4.13. Besides, the test accuracy of SGC surpass GCN on Citeseer. Thus, the reason why linear models sometimes perform is due to their smaller lipschitz continuity constants. (2) Compared with GCN, GCNII achieves smaller loss and

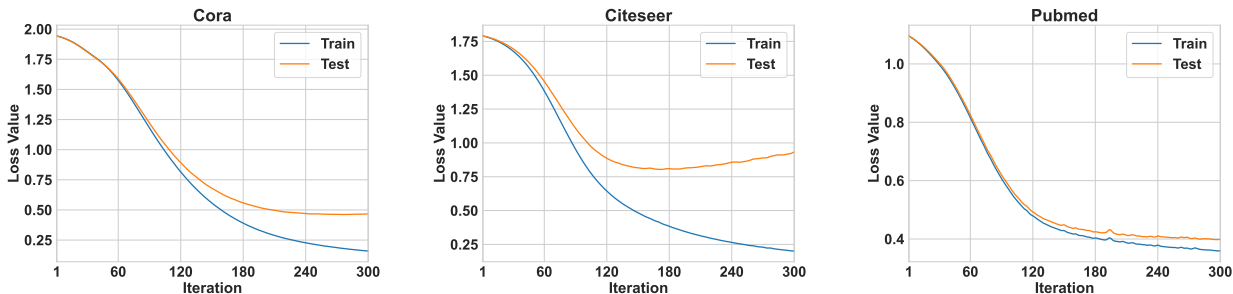


Figure 1. The loss value of GCN on training and test samples with the increase of iterations.

Table 1. Accuracy gap comparison of different baseline models on Cora, Citeseer and Pubmed.

	Cora	Citeseer	Pubmed
GCN	9.76±1.15	22.11±1.26	1.08±0.52
GCN*	13.45±1.28	26.48±1.21	1.49±0.63
GAT	11.00±0.75	22.69±0.84	1.52±0.43
GCNII	7.69±1.48	14.85±0.80	0.88±0.52
GCNII*	6.24±1.59	13.49±1.39	0.80±0.50
SGC	5.33±1.58	11.50±1.09	0.73±0.54
APPNP	7.72±1.54	9.99±1.17	0.85±0.46
GPR-GNN	8.90±1.22	19.08±0.95	0.96±0.49

Table 2. Test accuracy comparison of different baseline models on Cora, Citeseer and Pubmed.

	Cora	Citeseer	Pubmed
GCN	85.91±0.53	71.78±0.72	85.29±0.19
GCN*	82.49±0.59	66.74±1.07	84.21±0.26
GAT	86.10±0.51	72.90±0.65	85.45±0.26
GCNII	82.85±2.17	73.61±0.64	84.70±0.24
GCNII*	82.85±2.44	72.89±0.96	83.67±0.46
SGC	82.39±2.48	74.37±0.56	82.00±0.27
APPNP	79.14±3.17	74.12±0.62	82.86±0.29
GPR-GNN	87.24±0.71	73.79±0.67	85.07±0.34

accuracy gap with the same number of layers. We further estimate the generalization performance of GCN and GCNII with six layers (denoted as GCN\* and GCNII\*). Interestingly, with the increase of the number of hidden layers, the generalization performance of GCN decreases sharply. On the contrary, the loss and accuracy gap of GCNII remain unchanged. The test accuracy of GCNII also remain unchanged or only drops slightly. Therefore, the superior performance of GCNII comes from two perspectives: the first is learning non-degenerated representations by relieving over-smoothing and the second is robust generalization gap against the increase of the number of layers. (3) Although GPR-GNN achieve a competitive test accuracy, it has higher accuracy and loss gap than APPNP. Therefore,

Table 3. Loss gap comparison of different baseline models on Cora, Citeseer and Pubmed.

	Cora	Citeseer	Pubmed
GCN	0.30±0.03	0.77±0.04	0.03±0.01
GCN*	0.91±0.18	2.12±0.16	0.05±0.01
GAT	0.29±0.03	0.65±0.02	0.03±0.01
GCNII	0.19±0.03	0.43±0.02	0.02±0.01
GCNII*	0.16±0.03	0.43±0.03	0.02±0.01
SGC	0.12±0.03	0.28±0.02	0.01±0.00
APPNP	0.16±0.03	0.25±0.02	0.01±0.00
GPR-GNN	0.24±0.03	0.55±0.02	0.02±0.00

the unconstrained learning coefficients improve the fitting ability but also weaken generalization capability. Designing weight learning schema to balance the expressive and generalization could be a direction for spectral-based GNNs. (4) The generalization performance of GAT is slightly worse than GCN. Note that GAT is designed for inductive learning while our experimental setting is transductive. Thus, the superiority of GAT is not so obvious.

Besides, loss value of GCN on training and test samples w.r.t. iterations are presented in Figure 1. It can be seen that the loss gap increases with the increase of iterations, as demonstrated by Theorem 4.3. In general, the theoretical results are supported by the experimental results. It is worth pointing out that our analysis is only oriented to generalization gap. Smaller generalization gap does not necessarily mean better generalization ability, since the performance on test samples are determined by both training error and generalization gap.

## 6. Discussion and Conclusion

In this paper, we establish high probability learning guarantees for transductive SGD, by which the upper bound of generalization gap for some popular GNNs are derived. Experimental results on benchmark datasets support the theoretical results. This work sheds light on understanding the generalization of GNNs and provide some insights in



designing new GNN architecture with both expressiveness and generalization capabilities.

Although we have made efforts in generalization theory of GNNs, there are still some limitations in our analysis, which is left for future work to address: (1) The complexity based technique makes the dimension of parameters appearing in the bounds. Further research should focus on establishing dimension-independent bounds under milder assumption and deriving the lower bound that matches the upper bound. (2) We only analyze vanilla SGD in terms of optimization algorithms. Extending our results to SGD with momentum and adaptive learning rates is worth exploring. (3) Our analysis does not explicitly consider the heterophily of graphs. Deriving heterophily-dependent generalization bounds is an meaningful direction.

## References

- Bastings, J., Titov, I., Aziz, W., Marcheggiani, D., and Sima'an, K. Graph convolutional encoders for syntax-aware neural machine translation. In *Conference on Empirical Methods in Natural Language Processing*, pp. 1957–1967, 2017.
- Beck, D., Haffari, G., and Cohn, T. Graph-to-sequence learning using gated graph neural networks. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics*, pp. 273–283, 2018.
- Bevilacqua, B., Zhou, Y., and Ribeiro, B. Size-invariant graph representations for graph classification extrapolations. In *Proceedings of the 38th International Conference on Machine Learning*, pp. 837–851, 2021.
- Chen, M., Wei, Z., Huang, Z., Ding, B., and Li, Y. Simple and deep graph convolutional networks. In *Proceedings of the 37th International Conference on Machine Learning*, Proceedings of Machine Learning Research, pp. 1725–1735, 2020.
- Chen, Z., Villar, S., Chen, L., and Bruna, J. On the equivalence between graph isomorphism testing and function approximation with gnns. In *Advances in Neural Information Processing Systems*, pp. 15868–15876, 2019.
- Chien, E., Peng, J., Li, P., and Milenkovic, O. Adaptive universal generalized pagerank graph neural network. In *International Conference on Learning Representations*, 2021.
- Cong, W., Ramezani, M., and Mahdavi, M. On provable benefits of depth in training graph convolutional networks. In *Advances in Neural Information Processing Systems*, 2021.
- Defferrard, M., Bresson, X., and Vandergheynst, P. Convolutional neural networks on graphs with fast localized spectral filtering. In *Advances in Neural Information Processing Systems*, 2016.
- Dehmamy, N., Barabasi, A.-L., and Yu, R. Understanding the representation power of graph neural networks in learning graph topology. In *Advances in Neural Information Processing Systems*, 2019.
- Deng, L., Lian, D., Wu, C., and Chen, E. Graph convolution network based recommender systems: Learning guarantee and item mixture powered strategy. In *Advances in Neural Information Processing Systems*, 2022.
- Devroye, L., Györfi, L., and Lugosi, G. *A Probabilistic Theory of Pattern Recognition*. Springer, 1996.
- Du, S. S., Hou, K., Salakhutdinov, R., Póczos, B., Wang, R., and Xu, K. Graph neural tangent kernel: Fusing graph neural networks with graph kernels. In *Advances in Neural Information Processing Systems*, 2019.
- Duchi, J., Hazan, E., and Singer, Y. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of Machine Learning Research*, 12(61): 2121–2159, 2011.
- El-Yaniv, R. and Pechyony, D. Stable transductive learning. In *Conference on Learning Theory*, pp. 35–49, 2006.
- El-Yaniv, R. and Pechyony, D. Transductive rademacher complexity and its applications. In *Conference on Learning Theory*, pp. 157–171, 2007.
- Esser, P. M., Vankadara, L. C., and Ghoshdastidar, D. Learning theory can (sometimes) explain generalisation in graph neural networks. In *Advances in Neural Information Processing Systems*, pp. 27043–27056, 2021.
- Fan, W., Ma, Y., Li, Q., He, Y., Zhao, Y. E., Tang, J., and Yin, D. Graph neural networks for social recommendation. In *The World Wide Web Conference*, pp. 417–426, 2019.
- Fazlyab, M., Robey, A., Hassani, H., Morari, M., and Pappas, G. Efficient and accurate estimation of lipschitz constants for deep neural networks. In *Advances in Neural Information Processing Systems*, 2019.
- Federer, H. *Geometric Measure Theory*. Springer, 1969.
- Feng, J., Chen, Y., Li, F., Sarkar, A., and Zhang, M. How powerful are k-hop message passing graph neural networks. In *Advances in Neural Information Processing Systems*, 2022.
- Fey, M. and Lenssen, J. E. Fast graph representation learning with pytorch geometric. In *International Conference on Learning Representations*, 2019.

- Garg, V., Jegelka, S., and Jaakkola, T. Generalization and representational limits of graph neural networks. In *Proceedings of the 37th International Conference on Machine Learning*, pp. 3419–3430, 2020.
- Gasteiger, J., Bojchevski, A., and Günnemann, S. Predict then propagate: Graph neural networks meet personalized pagerank. In *International Conference on Learning Representations*, 2019.
- Giné, E. and Peña, V. H. *Decoupling: From Dependence to Independence*. Springer, 1999.
- Gori, M., Monfardini, G., and Scarselli, F. A new model for learning in graph domains. In *Proceedings. 2005 IEEE International Joint Conference on Neural Networks, 2005.*, pp. 729–734, 2005.
- Hamilton, W., Ying, Z., and Leskovec, J. Inductive representation learning on large graphs. In *Advances in Neural Information Processing Systems*, 2017.
- Han, J., Huang, W., Ma, H., Li, J., Tenenbaum, J. B., and Gan, C. Learning physical dynamics with subequivariant graph neural networks. In *Advances in Neural Information Processing Systems*, 2022.
- Hardt, M., Recht, B., and Singer, Y. Train faster, generalize better: Stability of stochastic gradient descent. In *Proceedings of the 33rd International Conference on Machine Learning*, pp. 1225–1234, 2016.
- He, M., Wei, Z., Huang, Z., and Xu, H. Bernnet: Learning arbitrary graph spectral filters via bernstein approximation. In *Advances in Neural Information Processing Systems*, 2021.
- He, X., Deng, K., Wang, X., Li, Y., Zhang, Y., and Wang, M. Lightgcn: Simplifying and powering graph convolution network for recommendation. In *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval*, pp. 639–648, 2020.
- Jacot, A., Gabriel, F., and Hongler, C. Neural tangent kernel: Convergence and generalization in neural networks. In *Advances in Neural Information Processing Systems*, 2018.
- Johnson, J., Gupta, A., and Fei-Fei, L. Image generation from scene graphs. In *Conference on Computer Vision and Pattern Recognition*, pp. 1219–1228, 2018.
- Ju, H., Li, D., Sharma, A., and Zhang, H. R. Generalization in graph neural networks: Improved pac-bayesian bounds on graph diffusion. In *Proceedings of The 26th International Conference on Artificial Intelligence and Statistics*, pp. 6314–6341, 2023.
- Keriven, N., Bietti, A., and Vaiter, S. Convergence and stability of graph convolutional networks on large random graphs. In *Advances in Neural Information Processing Systems*, 2020.
- Kingma, D. P. and Ba, J. Adam: A method for stochastic optimization. In *International Conference on Learning Representations*, 2015.
- Kipf, T. N. and Welling, M. Semi-supervised classification with graph convolutional networks. In *International Conference on Learning Representations*, 2017.
- Kuzborskij, I. and Lampert, C. Data-dependent stability of stochastic gradient descent. In *Proceedings of the 35th International Conference on Machine Learning*, pp. 2815–2824, 2018.
- Landrieu, L. and Simonovsky, M. Large-scale point cloud semantic segmentation with superpoint graphs. In *Conference on Computer Vision and Pattern Recognition*, pp. 4558–4567, 2018.
- Latała, R. and Oleszkiewicz, K. On the best constant in the khinchin-kahane inequality. *Studia Mathematica*, 109(1): 101–104, 1994.
- Lei, Y. and Tang, K. Learning rates for stochastic gradient descent with nonconvex objectives. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 43(12): 4505–4511, 2021.
- Li, H., Wang, M., Liu, S., Chen, P., and Xiong, J. Generalization guarantee of training graph convolutional networks with graph topology sampling. In *Proceedings of The 39th International Conference on Machine Learning*, pp. 13014–13051, 2022a.
- Li, H., Wang, X., Zhang, Z., and Zhu, W. Out-of-distribution generalization on graphs: A survey. *arXiv preprint arXiv:2202.07987*, 2022b.
- Li, Q., Han, Z., and Wu, X. Deeper insights into graph convolutional networks for semi-supervised learning. In *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence*, pp. 3538–3545, 2018.
- Li, S. and Liu, Y. Improved learning rates for stochastic optimization: Two theoretical viewpoints. *arXiv preprint arXiv:2107.08686*, 2021.
- Liao, R., Zhao, Z., Urtasun, R., and Zemel, R. Lanczosnet: Multi-scale deep graph convolutional networks. In *International Conference on Learning Representations*, 2019.
- Liao, R., Urtasun, R., and Zemel, R. A PAC-bayesian approach to generalization bounds for graph neural networks. In *International Conference on Learning Representations*, 2021.

- Liu, C., Zhu, L., and Belkin, M. Toward a theory of optimization for over-parameterized systems of non-linear equations: the lessons of deep learning. *arXiv preprint arXiv:2003.00307*, 2020.
- Luo, D., Cheng, W., Xu, D., Yu, W., Zong, B., Chen, H., and Zhang, X. Parameterized explainer for graph neural network. In *Advances in Neural Information Processing Systems*, 2020.
- Lv, S. Generalization bounds for graph convolutional neural networks via rademacher complexity. *arXiv preprint arXiv:2102.10234*, 2021.
- Maron, H., Ben-Hamu, H., Shamir, N., and Lipman, Y. Invariant and equivariant graph networks. In *International Conference on Learning Representations*, 2019.
- Oono, K. and Suzuki, T. Graph neural networks exponentially lose expressive power for node classification. In *International Conference on Learning Representations*, 2020a.
- Oono, K. and Suzuki, T. Optimization and generalization analysis of transduction through gradient boosting and application to multi-scale graph neural networks. In *Advances in Neural Information Processing Systems*, 2020b.
- Pfaff, T., Fortunato, M., Sanchez-Gonzalez, A., and Battaglia, P. W. Learning mesh-based simulation with graph networks. In *International Conference on Learning Representations*, 2021.
- Pisier, G. *The Volume of Convex Bodies and Banach Space Geometry*. Cambridge Tracts in Mathematics. Cambridge University Press, 1989.
- Qi, X., Liao, R., Jia, J., Fidler, S., and Urtasun, R. 3d graph neural networks for RGBD semantic segmentation. In *International Conference on Computer Vision*, pp. 5209–5218, 2017.
- Rong, Y., Huang, W., Xu, T., and Huang, J. Dropedge: Towards deep graph convolutional networks on node classification. In *International Conference on Learning Representations*, 2020.
- Sanchez-Gonzalez, A., Godwin, J., Pfaff, T., Ying, R., Leskovec, J., and Battaglia, P. W. Learning to simulate complex physics with graph networks. In *Proceedings of the 37th International Conference on Machine Learning*, volume 119, pp. 8459–8468, 2020.
- Satorras, V. G. and Estrach, J. B. Few-shot learning with graph neural networks. In *International Conference on Learning Representations*, 2018.
- Scarselli, F., Gori, M., Tsoi, A. C., Hagenbuchner, M., and Monfardini, G. The graph neural network model. *IEEE Transactions on Neural Networks*, 20(1):61–80, 2009.
- Scarselli, F., Tsoi, A. C., and Hagenbuchner, M. The vapnik–chervonenkis dimension of graph and recursive neural networks. *Neural Networks*, 108:248–259, 2018.
- Sen, P., Namata, G., Bilgic, M., Getoor, L., Gallagher, B., and Eliassi-Rad, T. Collective classification in network data. *AI magazine*, 29(3):93–106, 2008.
- Shen, Z.-A., Luo, T., Zhou, Y.-K., Yu, H., and Du, P.-F. Npi-gnn: Predicting ncna-protein interactions with deep graph neural networks. *Briefings in bioinformatics*, 2021.
- Song, L., Zhang, Y., Wang, Z., and Gildea, D. A graph-to-sequence model for amr-to-text generation. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics*, pp. 1616–1626, 2018.
- Vapnik, V. *Statistical learning theory*. Wiley, 1998.
- Vapnik, V. *Estimation of Dependences Based on Empirical Data, Second Edition*. Springer, 2006.
- Veličković, P., Cucurull, G., Casanova, A., Romero, A., Liò, P., and Bengio, Y. Graph attention networks. In *International Conference on Learning Representations*, 2018.
- Verma, S. and Zhang, Z.-L. Stability and generalization of graph convolutional neural networks. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pp. 1539–1548, 2019.
- Virmaux, A. and Scaman, K. Lipschitz regularity of deep neural networks: analysis and efficient estimation. In *Advances in Neural Information Processing Systems*, 2018.
- Vu, M. and Thai, M. T. Pgm-explainer: Probabilistic graphical model explanations for graph neural networks. In *Advances in Neural Information Processing Systems*, pp. 12225–12235, 2020.
- Wang, Y., Wang, Y., Yang, J., and Lin, Z. Dissecting the diffusion process in linear graph convolutional networks. In *Advances in Neural Information Processing Systems*, 2021.
- Wu, F., Souza, A., Zhang, T., Fifty, C., Yu, T., and Weinberger, K. Simplifying graph convolutional networks. In *Proceedings of the 36th International Conference on Machine Learning*, pp. 6861–6871, 2019.
- Wu, Q., Zhang, H., Yan, J., and Wipf, D. Handling distribution shifts on graphs: An invariance perspective. In *International Conference on Learning Representations*, 2022.

- Wu, Q., Chen, Y., Yang, C., and Yan, J. Energy-based out-of-distribution detection for graph neural networks. In *The Eleventh International Conference on Learning Representations*, 2023.
- Wu, Z., Pan, S., Chen, F., Long, G., Zhang, C., and Yu, P. S. A comprehensive survey on graph neural networks. *IEEE Transactions on Neural Networks and Learning Systems*, 32(1):4–24, 2021.
- Xu, K., Li, C., Tian, Y., Sonobe, T., Kawarabayashi, K., and Jegelka, S. Representation learning on graphs with jumping knowledge networks. In *Proceedings of the 35th International Conference on Machine Learning*, volume 80, pp. 5449–5458, 2018.
- Xu, K., Hu, W., Leskovec, J., and Jegelka, S. How powerful are graph neural networks? In *International Conference on Learning Representations*, 2019.
- Xu, K., Li, J., Zhang, M., Du, S. S., ichi Kawarabayashi, K., and Jegelka, S. What can neural networks reason about? In *International Conference on Learning Representations*, 2020.
- Xu, K., Zhang, M., Li, J., Du, S. S., Kawarabayashi, K.-I., and Jegelka, S. How neural networks extrapolate: From feedforward to graph neural networks. In *International Conference on Learning Representations*, 2021.
- Xu, Y. and Zeevi, A. Towards optimal problem dependent generalization error bounds in statistical learning theory. *arXiv preprint arXiv:2011.06186*, 2020.
- Yang, C., Wu, Q., Wang, J., and Yan, J. Graph neural networks are inherently good generalizers: Insights by bridging GNNs and MLPs. In *The Eleventh International Conference on Learning Representations*, 2023.
- Yang, N., Zeng, K., Wu, Q., Jia, X., and Yan, J. Learning substructure invariance for out-of-distribution molecular representations. In *Advances in Neural Information Processing Systems*, 2022.
- Yang, Z., Cohen, W. W., and Salakhutdinov, R. Revisiting semi-supervised learning with graph embeddings. In *Proceedings of the 33rd International Conference on Machine Learning*, pp. 40–48, 2016.
- Yehudai, G., Fetaya, E., Meiri, E. A., Chechik, G., and Maron, H. From local structures to size generalization in graph neural networks. In *Proceedings of the 38th International Conference on Machine Learning*, pp. 11975–11986, 2021.
- Ying, R., He, R., Chen, K., Eksombatchai, P., Hamilton, W. L., and Leskovec, J. Graph convolutional neural networks for web-scale recommender systems. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pp. 974–983, 2018.
- Ying, Y. and Campbell, C. Rademacher chaos complexities for learning the kernel problem. *Neural Computation*, 22(11):2858–2886, 2010.
- Ying, Z., Bourgeois, D., You, J., Zitnik, M., and Leskovec, J. Gnnexplainer: Generating explanations for graph neural networks. In *Advances in Neural Information Processing Systems*, pp. 9240–9251, 2019.
- Yuan, H., Tang, J., Hu, X., and Ji, S. Xggn: Towards model-level explanations of graph neural networks. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pp. 430–438, 2020.
- Yuan, H., Yu, H., Wang, J., Li, K., and Ji, S. On explainability of graph neural networks via subgraph explorations. In *Proceedings of the 38th International Conference on Machine Learning*, pp. 12241–12252, 2021.
- Zeng, H., Zhou, H., Srivastava, A., Kannan, R., and Prasanna, V. Graphsaint: Graph sampling based inductive learning method. In *International Conference on Learning Representations*, 2020.
- Zeng, H., Zhang, M., Xia, Y., Srivastava, A., Malevich, A., Kannan, R., Prasanna, V. K., Jin, L., and Chen, R. Decoupling the depth and scope of graph neural networks. In *Advances in Neural Information Processing Systems*, pp. 19665–19679, 2021.
- Zhang, S., Wang, M., Liu, S., Chen, P.-Y., and Xiong, J. Fast learning of graph neural networks with guaranteed generalizability: One-hidden-layer case. In *Proceedings of the 37th International Conference on Machine Learning*, pp. 11268–11277, 2020.
- Zhang, Z., Cui, P., and Zhu, W. Deep learning on graphs: A survey. *IEEE Transactions on Knowledge and Data Engineering*, 34(1):249–270, 2022.
- Zhao, L. and Akoglu, L. Pairnorm: Tackling oversmoothing in gnns. In *International Conference on Learning Representations*, 2020.
- Zhou, D., Tang, Y., Yang, Z., Cao, Y., and Gu, Q. On the convergence of adaptive gradient methods for nonconvex optimization. *arXiv preprint arXiv:1808.05671*, 2018.
- Zhou, J., Cui, G., Hu, S., Zhang, Z., Yang, C., Liu, Z., Wang, L., Li, C., and Sun, M. Graph neural networks: A review of methods and applications. *AI Open*, 1:57–81, 2020.
- Zhou, X. and Wang, H. The generalization error of graph convolutional networks may enlarge with more layers. *Neurocomputing*, 424:97–106, 2021.

Zhu, H. and Koniusz, P. Simple spectral graph convolution. In *International Conference on Learning Representations*, 2021.

Zhu, Y., Xu, W., Zhang, J., Du, Y., Zhang, J., Liu, Q., Yang, C., and Wu, S. A survey on graph structure learning: Progress and opportunities. *arXiv preprint arXiv:2103.03036*, 2021.

## A. Notations and Lemmas

In this section, we will present some notations, definitions and lemmas that will be used in subsequent analysis. Let  $f : \mathbb{R}^{m \times n} \mapsto \mathbb{R}$  be a real-value function with variable  $\mathbf{W} \in \mathbb{R}^{m \times n}$ . We stipulate that

$$\nabla_{\text{vec}[\mathbf{W}]} f = \left[ \frac{\partial f}{\partial W_{11}}, \dots, \frac{\partial f}{\partial W_{m1}}, \dots, \frac{\partial f}{\partial W_{1n}}, \dots, \frac{\partial f}{\partial W_{mn}} \right]^\top \in \mathbb{R}^{mn \times 1}.$$

Denote by  $\frac{\partial f}{\partial \text{vec}[\mathbf{W}]}$  the Jacobian matrix, we have  $\frac{\partial f}{\partial \text{vec}[\mathbf{W}]} = \nabla_{\text{vec}[\mathbf{W}]}^\top f \in \mathbb{R}^{1 \times mn}$ . Denote by  $\mathbf{W}_{i*}$  the  $i$ -th row of matrix  $\mathbf{W}$ . We use  $\odot$  and  $\otimes$  to denote Hadamard product and Kronecker product, respectively. The activation function  $\sigma(\cdot)$  in this work is defined as

$$\sigma(x) = \begin{cases} 0, & x \leq 0, \\ x^q, & 0 < x \leq \left(\frac{1}{q}\right)^{\frac{1}{q-1}}, \\ x - \left(\frac{1}{q}\right)^{\frac{1}{q-1}} + \left(\frac{1}{q}\right)^{\frac{q}{q-1}}, & x > \left(\frac{1}{q}\right)^{\frac{1}{q-1}}, \end{cases}$$

where  $q \in (1, 2]$ . It can be verify that this activation is differential on  $\mathbb{R}$ , and its derivation is

$$\sigma'(x) = \begin{cases} 0, & x \leq 0, \\ qx^{q-1}, & 0 < x \leq \left(\frac{1}{q}\right)^{\frac{1}{q-1}}, \\ 1, & x > \left(\frac{1}{q}\right)^{\frac{1}{q-1}}. \end{cases}$$

When setting  $p \approx 1$  (e.g.,  $q = 1.1$ ), this activation function has tolerate approximation error to vanilla ReLU function. Now we show some property of  $\sigma(\cdot)$  that used in the sequential proofs.

- $\|\sigma(\mathbf{u})\|_2 < \|\mathbf{u}\|_2$  for any  $\mathbf{u} \in \mathbb{R}^d$ . We only need to show that  $|\sigma(u_i)| \leq |u_i|$  holds for  $i \in [d]$ . The case that  $u_i \in (-\infty, 0]$  is trivial. If  $u_i \in (0, (1/q)^{1/(q-1)}]$ , since  $q > 1$  and  $(1/q)^{1/(q-1)} < 1$ , we have  $|\sigma(u_i)| = u_i^q \leq u_i = |u_i|$ . If  $u_i \in ((1/q)^{1/(q-1)}, \infty)$ , note that  $(1/q)^{q/(q-1)} < (1/q)^{1/(q-1)}$ , we have  $|\sigma(u_i)| \leq u_i = |u_i|$ .
- $\|\sigma'(\mathbf{u}) \odot \mathbf{v}\|_2 \leq \|\mathbf{v}\|_2$  for any  $\mathbf{u}, \mathbf{v} \in \mathbb{R}^d$ . By the formulation of  $\sigma'(x)$ , we have  $|\sigma'(x)| \leq 1$ . Then

$$\|\sigma'(\mathbf{u}) \odot \mathbf{v}\|_2 = \sqrt{\sum_{i=1}^d |\sigma'(u_i)|^2 |v_i|^2} \leq \sqrt{\sum_{i=1}^d |v_i|^2} = \|\mathbf{v}\|_2.$$

- $\|\sigma'(\mathbf{u}) - \sigma'(\mathbf{v})\|_2 \leq qd^{\frac{2-q}{2}} \|\mathbf{u} - \mathbf{v}\|_2^{q-1}$  for any  $\mathbf{u}, \mathbf{v} \in \mathbb{R}^d$ . We first show that for any  $x, y \in \mathbb{R}$ ,  $|\sigma'(x) - \sigma'(y)| \leq q|x - y|^{q-1}$  holds. The case that  $x, y \in (-\infty, 0]$  and  $x, y \in [(1/q)^{1/(q-1)}, \infty)$  are trivial. If  $x, y \in (0, (1/q)^{1/(q-1)}]$ , we have  $|qx^{q-1} - qy^{q-1}| \leq q|x - y|^{q-1}$ . If  $x \in (-\infty, 0]$  and  $y \in (0, (1/q)^{1/(q-1)}]$ , we have  $|qy^{q-1}| = qy^{q-1} \leq q|x - y|^{q-1}$ . If  $x \in (0, (1/q)^{1/(q-1)}]$  and  $y \in ((1/q)^{1/(q-1)}, \infty)$ , we have  $|qx^{q-1} - 1| \leq |qx^{q-1} - qy^{q-1}| \leq q|x - y|^{q-1}$ . If  $x \in (-\infty, 0]$  and  $y \in ((1/q)^{1/(q-1)}, \infty)$ , we have  $|\sigma'(x) - \sigma'(y)| \leq qy^{q-1} \leq q|x - y|^{q-1}$ . Thus,

$$\begin{aligned} \|\sigma'(\mathbf{u}) - \sigma'(\mathbf{v})\|_2 &= \sqrt{\sum_{i=1}^d |\sigma'(u_i) - \sigma'(v_i)|^2} \leq \sqrt{\sum_{i=1}^d q^2 |u_i - v_i|^{2(q-1)}} \\ &\leq \sqrt{\left(\sum_{i=1}^d |u_i - v_i|^2\right)^{q-1} \left(\sum_{i=1}^d q^{\frac{2}{2-q}}\right)^{2-q}} \leq qd^{\frac{2-q}{2}} \|\mathbf{u} - \mathbf{v}\|_2^{q-1}. \end{aligned}$$

With the above notations, we give the following lemmas.

**Lemma A.1.** Denote by  $\tilde{\mathbf{A}}$  the normalized adjacency matrix with self-loop, we have  $\|\tilde{\mathbf{A}}\|_\infty \leq \sqrt{\frac{\text{deg}_{\max} + 1}{\text{deg}_{\min} + 1}}$ .

*Proof.* By definition,  $\tilde{\mathbf{A}}_{ij} \geq 0$  holds for any  $i, j \in [n]$ . For any fixed  $i \in [n]$ , let  $\mathcal{N}_i$  be the index set of the  $i$ -th nodes' one-hop neighbors, we have

$$\begin{aligned} \sum_{j=1}^n \tilde{\mathbf{A}}_{ij} &= \sum_{j=1}^n \frac{\mathbf{A}_{ij}}{\sqrt{\deg_i + 1} \sqrt{\deg_j + 1}} \\ &= \frac{1}{\sqrt{\deg_i + 1}} \left( \frac{1}{\sqrt{\deg_i + 1}} + \sum_{j \in \mathcal{N}_i} \frac{1}{\sqrt{\deg_j + 1}} \right) \\ &\leq \frac{1}{\sqrt{\deg_i + 1}} \left( \frac{1}{\sqrt{\deg_{\min} + 1}} + \sum_{j \in \mathcal{N}_i} \frac{1}{\sqrt{\deg_{\min} + 1}} \right) \\ &\leq \frac{1}{\sqrt{\deg_i + 1}} \frac{\deg_i + 1}{\sqrt{\deg_{\min} + 1}} = \frac{\sqrt{\deg_i + 1}}{\sqrt{\deg_{\min} + 1}} \leq \sqrt{\frac{\deg_{\max} + 1}{\deg_{\min} + 1}}. \end{aligned}$$

**Lemma A.2.** Denote by  $\mathbf{u} \in \mathbb{R}^m$ ,  $\mathbf{v} \in \mathbb{R}^n$ , we have  $\|\mathbf{u} \otimes \mathbf{v}\|_2 = \|\mathbf{u}\|_2 \|\mathbf{v}\|_2$ .

*Proof.* One can find that

$$\|\mathbf{u} \otimes \mathbf{v}\|_2 = \sqrt{\sum_{j=1}^m \|u_j \mathbf{v}\|_2^2} = \sqrt{\sum_{j=1}^m u_j^2 \|\mathbf{v}\|_2^2} = \sqrt{\|\mathbf{u}\|_2^2 \|\mathbf{v}\|_2^2} = \|\mathbf{u}\|_2 \|\mathbf{v}\|_2.$$

**Lemma A.3.** Denote by  $\mathbf{W}_1, \mathbf{W}_2 \in \mathbb{R}^{m \times n}$ , we have

$$\|\mathbf{W}_1 - \mathbf{W}_2\| \leq \|\text{vec}[\mathbf{W}_1] - \text{vec}[\mathbf{W}_2]\|_2.$$

*Proof.* We have

$$\begin{aligned} \|\mathbf{W}_1 - \mathbf{W}_2\| &= \sup_{\|\mathbf{u}\|_2=1} \|(\mathbf{W}_1 - \mathbf{W}_2)\mathbf{u}\|_2 \triangleq \|(\mathbf{W}_1 - \mathbf{W}_2)\mathbf{u}^*\|_2 \\ &= \sqrt{\sum_{j=1}^m ([\mathbf{W}_1]_{j,:} \mathbf{u}^* - [\mathbf{W}_2]_{j,:} \mathbf{u}^*)^2} \leq \sqrt{\sum_{j=1}^m \|[\mathbf{W}_1]_{j,:} - [\mathbf{W}_2]_{j,:}\|_2^2} = \|\text{vec}(\mathbf{W}_1) - \text{vec}(\mathbf{W}_2)\|_2, \end{aligned}$$

where the last inequality follows from the Cauchy-Schwarz inequality and  $\|\mathbf{u}^*\|_2 = 1$ . This finishes the proof.

**Lemma A.4.** Denote by  $\{\mathbf{W}_h\}_{h=1}^H$  the learnable parameters (w.l.o.g. we assume that each parameter is matrix since vector is a special case of matrix). If for  $h \in [H]$ ,

$$|\ell(\mathbf{W}_1, \dots, \mathbf{W}_h, \dots, \mathbf{W}_H) - \ell(\mathbf{W}_1, \dots, \mathbf{W}'_h, \dots, \mathbf{W}_H)| \leq L_h \|\text{vec}[\mathbf{W}_h] - \text{vec}[\mathbf{W}'_h]\|_2,$$

and

$$\left\| \frac{\partial \ell(\mathbf{W}_1, \dots, \mathbf{W}_H)}{\partial \text{vec}[\mathbf{W}_h]} - \frac{\partial \ell(\mathbf{W}'_1, \dots, \mathbf{W}'_H)}{\partial \text{vec}[\mathbf{W}_h]} \right\|_2 \leq \sum_{i=1}^H \left[ P_{hi} \|\text{vec}[\mathbf{W}_i] - \text{vec}[\mathbf{W}'_i]\|_2 + \tilde{P}_{hi} \|\text{vec}[\mathbf{W}_i] - \text{vec}[\mathbf{W}'_i]\|_2^{\tilde{\alpha}} \right],$$

then there exist  $P, L > 0$  such that  $|\ell(\mathbf{w}) - \ell(\mathbf{w}')| \leq L \|\mathbf{w} - \mathbf{w}'\|_2$  and  $\|\nabla \ell(\mathbf{w}) - \nabla \ell(\mathbf{w}')\|_2 \leq P_{\mathcal{F}} \max\{\|\mathbf{w} - \mathbf{w}'\|_2, \|\mathbf{w} - \mathbf{w}'\|_2^{\tilde{\alpha}}\}$  hold, where  $\mathbf{w} = [\text{vec}[\mathbf{W}_1]; \dots; \text{vec}[\mathbf{W}_H]]$ .

*Proof.* By definition, the gradient of  $\ell$  w.r.t  $\mathbf{w}$  is  $\nabla\ell(\mathbf{w}) = \left[ \frac{\partial\ell}{\partial\text{vec}[\mathbf{W}_1]}, \dots, \frac{\partial\ell}{\partial\text{vec}[\mathbf{W}_H]} \right]^\top$ . Then we have

$$\begin{aligned} |\ell(\mathbf{w}) - \ell(\mathbf{w}')| &\leq \sum_{h=1}^H |\ell(\mathbf{W}_1, \dots, \mathbf{W}_h, \dots, \mathbf{W}_H) - \ell(\mathbf{W}_1, \dots, \mathbf{W}'_h, \dots, \mathbf{W}_H)| \\ &\leq \sum_{h=1}^H L_h \left\| \text{vec}[\mathbf{W}_h] - \text{vec}[\tilde{\mathbf{W}}_h] \right\|_2 \\ &\leq \left( \sum_{h=1}^H L_h^2 \right)^{\frac{1}{2}} \left( \sum_{h=1}^H \left\| \text{vec}[\mathbf{W}_h] - \text{vec}[\mathbf{W}'_h] \right\|_2^2 \right)^{\frac{1}{2}} \\ &= L \|\mathbf{w} - \tilde{\mathbf{w}}\|_2, \end{aligned}$$

where we obtain the last inequality by Cauchy-Schwarz inequality. Similarly, we have

$$\begin{aligned} &\|\nabla\ell(\mathbf{w}) - \nabla\ell(\mathbf{w}')\|_2 \\ &\leq \left\| \frac{\partial\ell(\mathbf{W}_1, \dots, \mathbf{W}_H)}{\partial\text{vec}[\mathbf{W}_h]} - \frac{\partial\ell(\mathbf{W}'_1, \dots, \mathbf{W}'_H)}{\partial\text{vec}[\mathbf{W}_h]} \right\|_2 \\ &\leq \sum_{h=1}^H \left[ \sum_{i=1}^H P_{hi} \left\| \text{vec}[\mathbf{W}_i] - \text{vec}[\mathbf{W}'_i] \right\|_2 \right] + \sum_{h=1}^H \left[ \sum_{i=1}^H \tilde{P}_{hi} \left\| \text{vec}[\mathbf{W}_i] - \text{vec}[\mathbf{W}'_i] \right\|_2^{\tilde{\alpha}} \right] \\ &= \sum_{i=1}^H \left[ \sum_{h=1}^H P_{hi} \left\| \text{vec}[\mathbf{W}_i] - \text{vec}[\mathbf{W}'_i] \right\|_2 \right] + \sum_{i=1}^H \left[ \sum_{h=1}^H \tilde{P}_{hi} \left\| \text{vec}[\mathbf{W}_i] - \text{vec}[\mathbf{W}'_i] \right\|_2^{\tilde{\alpha}} \right] \\ &= \sum_{i=1}^H P_i \left\| \text{vec}[\mathbf{W}_i] - \text{vec}[\mathbf{W}'_i] \right\|_2 + \sum_{i=1}^H \tilde{P}_i \left\| \text{vec}[\mathbf{W}_i] - \text{vec}[\mathbf{W}'_i] \right\|_2^{\tilde{\alpha}} \\ &\leq \left( \sum_{i=1}^H P_i^2 \right)^{\frac{1}{2}} \left( \sum_{i=1}^H \left\| \text{vec}[\mathbf{W}_i] - \text{vec}[\mathbf{W}'_i] \right\|_2^2 \right)^{\frac{1}{2}} + \left( \sum_{i=1}^H P_i^{2-\frac{2}{\tilde{\alpha}}} \right)^{1-\frac{\tilde{\alpha}}{2}} \left( \sum_{i=1}^H \left\| \text{vec}[\mathbf{W}_i] - \text{vec}[\mathbf{W}'_i] \right\|_2^2 \right)^{\frac{\tilde{\alpha}}{2}} \\ &= \left( \sum_{i=1}^H P_i^2 \right)^{\frac{1}{2}} \|\mathbf{w} - \mathbf{w}'\|_2 + \left( \sum_{i=1}^H P_i^{2-\frac{2}{\tilde{\alpha}}} \right)^{1-\frac{\tilde{\alpha}}{2}} \|\mathbf{w} - \mathbf{w}'\|_2^{\tilde{\alpha}} \\ &\leq \max \left\{ \left( \sum_{i=1}^H P_i^2 \right)^{\frac{1}{2}} + \left( \sum_{i=1}^H P_i^{2-\frac{2}{\tilde{\alpha}}} \right)^{1-\frac{\tilde{\alpha}}{2}} \right\} \max \left\{ \|\mathbf{w} - \mathbf{w}'\|_2, \|\mathbf{w} - \mathbf{w}'\|_2^{\tilde{\alpha}} \right\} \\ &= P \max \left\{ \|\mathbf{w} - \mathbf{w}'\|_2, \|\mathbf{w} - \mathbf{w}'\|_2^{\tilde{\alpha}} \right\}. \end{aligned} \tag{10}$$

where we define  $P_i = \sum_{h=1}^H P_{hi}$  and  $\tilde{P}_i = \sum_{h=1}^H \tilde{P}_{hi}$ . The second inequality is due to the Hölder inequality.

**Lemma A.5.** Denote by  $\mathbf{v} \in \mathbb{R}^d$ . Let  $f: \mathbb{R}^d \mapsto \mathbb{R}^d$  be  $f(\mathbf{v})_j = \frac{e^{v_j}}{\sum_{i=1}^d e^{v_i}}$ . For any  $\mathbf{v}, \mathbf{v}' \in \mathbb{R}^d$ , we have  $\|f(\mathbf{v}) - f(\mathbf{v}')\|_2 \leq 2\|\mathbf{v} - \mathbf{v}'\|_2$ .

*Proof.* By (Federer, 1969), we have  $\|f(\mathbf{v}) - f(\mathbf{v}')\|_2 \leq \sup_{\mathbf{v} \in \mathbb{R}^d} \|J(\mathbf{v})\| \|\mathbf{v} - \mathbf{v}'\|_2$ , where  $J$  is the Jacobian. For the aforementioned  $f$ , we have  $J(\mathbf{v}) = \text{diag}(f(\mathbf{v})) - f(\mathbf{v})f(\mathbf{v})^\top$ . Then

$$\|J(\mathbf{v})\| = \|\text{diag}(f(\mathbf{v})) - f(\mathbf{v})f(\mathbf{v})^\top\| \leq \|\text{diag}(f(\mathbf{v}))\| + \|f(\mathbf{v})f(\mathbf{v})^\top\|. \tag{11}$$

First,

$$\|\text{diag}(f(\mathbf{v}))\| = \sup_{\|\mathbf{w}\|_2=1} \|\text{diag}(f(\mathbf{v})) \mathbf{w}\|_2 = \sup_{\|\mathbf{w}\|_2=1} \sqrt{\sum_{i=1}^d f^2(\mathbf{v})_i w_i^2} = \max_{i \in [d]} f(\mathbf{v})_i \leq 1. \tag{12}$$

Besides,

$$\|f(\mathbf{v})f(\mathbf{v})^\top\| = \sup_{\|\mathbf{w}\|_2=1} \|f(\mathbf{v})f(\mathbf{v})^\top \mathbf{w}\|_2 = \|f(\mathbf{v})\|_2 \sup_{\|\mathbf{w}\|_2=1} |f(\mathbf{v})^\top \mathbf{w}| \leq \|f(\mathbf{v})\|_2^2 \leq 1, \tag{13}$$



where the last inequality is due to  $\sum_{i=1}^d f(\mathbf{v})_i = 1$ . Plugging Eq. (12) and Eq. (13) into Eq. (11), the proof is completed.

**Lemma A.6.** Denote by  $\mathbf{v} \in \mathbb{R}^d$ . Let  $f : \mathbb{R}^d \mapsto \mathbb{R}$  be  $f(\mathbf{v})_j = -\sum_{k=1}^K y_k \log \hat{y}_k$ , where  $\hat{y}_k = \frac{e^{v_k}}{\sum_{i=1}^K e^{v_i}}$ . For any  $\mathbf{v}, \mathbf{v}' \in \mathbb{R}^d$ , we have  $|f(\mathbf{v}) - f(\mathbf{v}')| \leq \sqrt{2} \|\mathbf{v} - \mathbf{v}'\|_2$ .

*Proof.* By the chain rule, the Jacobian is  $J(\mathbf{v}) = \hat{\mathbf{y}}(\mathbf{v}) - \mathbf{y}$ . Note that  $|f(\mathbf{v}) - f(\mathbf{v}')| \leq \sup_{\mathbf{v} \in \mathbb{R}^d} \|J(\mathbf{v})\|_2 \|\mathbf{v} - \mathbf{v}'\|_2$ . W.o.l.g, we assume that  $y_1 = 1$ , then

$$\|J(\mathbf{v})\|_2 = \sqrt{\sum_{k=1}^K (\hat{y}_k - y_k)^2} = \sqrt{(1 - \hat{y}_1)^2 + \sum_{k=2}^K \hat{y}_k^2} \leq \sqrt{1 + \sum_{k=1}^K \hat{y}_k^2} \leq \sqrt{2},$$

where the last inequality is due to  $\sum_{k=1}^K \hat{y}_k = 1$ .

**Lemma A.7.** Let  $\mathcal{F} : \mathcal{W} \times \mathcal{Z} \times \mathcal{Z} \mapsto \mathbb{R}$  be a parametric function class. For  $\mathbf{w}, \mathbf{w}' \in \mathcal{W}$ , the empirical metric defined on  $\mathcal{F}$  is defined as

$$d_S(\mathbf{w}, \mathbf{w}') = \left( \frac{1}{n^2} \sum_{1 \leq i < j \leq n} |f(\mathbf{w}; z_i, z_j) - f(\mathbf{w}'; z_i, z_j)|^2 \right)^{\frac{1}{2}}.$$

For specific  $\mathbf{w}^{(1)} \in \mathcal{W}$ , assume that  $\sup_{\mathbf{w} \in \mathcal{W}} d_S(\mathbf{w}, \mathbf{w}^{(1)}) \leq D$  and  $|f(\mathbf{w}^{(1)}; z_i, z_j)| \leq M_0$  hold. Then we have

$$\mathcal{U}(\mathcal{F}) \triangleq \frac{1}{n} \mathbb{E}_{\boldsymbol{\sigma}} \left[ \sup_{\mathbf{w} \in \mathcal{W}} \sum_{1 \leq i < j \leq n} \sigma_i \sigma_j f(\mathbf{w}; z_i, z_j) \right] \leq M_0 + 24e \int_0^D (1 + \log \mathcal{N}(r, \mathcal{F}_{\mathcal{W}}, d_S)) dr,$$

where  $\boldsymbol{\sigma}$  is the transductive Rademacher variable.

*Proof.* The proof extend Theorem 2 in (Ying & Campbell, 2010) to the transductive Rademacher chaos complexity. The first step is to show that the following inequality holds for  $1 < p \leq q < \infty$ ,  $d \geq 1$  and  $\gamma = \left(\frac{p-1}{q-1}\right)^{\frac{1}{2}}$ :

$$\begin{aligned} & \left[ \mathbb{E} \left\| x + \gamma \sum_{i=1}^n x_i \sigma_i + \gamma^2 \sum_{i_1 < i_2 \leq n} x_{i_1 i_2} \sigma_{i_1} \sigma_{i_2} + \cdots + \gamma^d \sum_{i_1 < \cdots < i_d \leq n} x_{i_1 \dots i_d} \sigma_{i_1} \cdots \sigma_{i_d} \right\|_2^q \right]^{\frac{1}{q}} \\ & \leq \left[ \mathbb{E} \left\| x + \sum_{i=1}^n x_i \epsilon_i + \sum_{i_1 < i_2 \leq n} x_{i_1 i_2} \epsilon_{i_1} \epsilon_{i_2} + \cdots + \sum_{i_1 < \cdots < i_d \leq n} x_{i_1 \dots i_d} \epsilon_{i_1} \cdots \epsilon_{i_d} \right\|_2^p \right]^{\frac{1}{p}}, \end{aligned} \quad (14)$$

where  $\sigma$  and  $\epsilon$  are transductive and standard Rademacher variable, respectively. The process generally follows that of Theorem 3.2.2 in (Giné & Peña, 1999). First, consider the case that  $n = 1$ , we have to show that  $(\mathbb{E}|x + \gamma \sigma y|^q)^{\frac{1}{q}} \leq (\mathbb{E}|x + \epsilon y|^p)^{\frac{1}{p}}$  holds. This inequality naturally holds when  $x = y = 0$  or  $y = 0$ . When  $x = 0$  and  $y \neq 0$ , we have

$$(\mathbb{E}|x + \gamma \sigma y|^q)^{\frac{1}{q}} \Big|_{x=0} = (2p_0 |\gamma y|^q)^{\frac{1}{q}} \leq (|\gamma y|^q)^{\frac{1}{q}} \leq |\gamma y| = (\mathbb{E}|x + \gamma \epsilon y|^q)^{\frac{1}{q}} \Big|_{x=0},$$

where the inequality is due to  $p_0 \leq \frac{1}{2}$ . When  $x \neq 0$  and  $y \neq 0$ , let  $u = \frac{y}{x}$ , then

$$(\mathbb{E}|x + \gamma \sigma y|^q)^{\frac{1}{q}} \leq (\mathbb{E}|x + \epsilon y|^p)^{\frac{1}{p}} \Leftrightarrow (\mathbb{E}|1 + \gamma \sigma u|^q)^{\frac{1}{q}} \leq (\mathbb{E}|1 + \gamma \sigma u|^p)^{\frac{1}{p}}.$$

By symmetric, we only have to discuss the case that  $u \geq 0$ . For  $0 \leq u \leq 1$ :

$$\begin{aligned}
 (\mathbb{E}|1 + \gamma\sigma u|^q)^{\frac{1}{q}} &= (p_0|1 + \gamma u|^q + p_0|1 - \gamma u|^q + (1 - 2p_0))^{\frac{1}{q}} \\
 &= \left[ p_0 + \sum_{k=1}^{\infty} p_0 \binom{q}{k} \gamma^k u^k + p_0 + \sum_{k=1}^{\infty} p_0 \binom{q}{k} (-1)^k \gamma^k u^k + (1 - 2p_0) \right]^{\frac{1}{q}} \\
 &= \left[ 2p_0 + 2p_0 \sum_{k=1}^{\infty} \binom{q}{2k} \gamma^{2k} u^{2k} + (1 - 2p_0) \right]^{\frac{1}{q}} \\
 &= \left[ 1 + 2p_0 \sum_{k=1}^{\infty} \binom{q}{2k} \gamma^{2k} u^{2k} \right]^{\frac{1}{q}} \leq \left[ 1 + \sum_{k=1}^{\infty} \binom{q}{2k} \gamma^{2k} u^{2k} \right]^{\frac{1}{q}} \\
 &= \left[ \frac{1}{2} |1 + \gamma u|^q + |1 - \gamma u|^q \right]^{\frac{1}{q}} \leq (\mathbb{E}|1 + \epsilon u|^q)^{\frac{1}{q}},
 \end{aligned}$$

where the first inequality is due to  $p_0 \leq \frac{1}{2}$ , and the last inequality is from Eq. (3.2.4') in (Giné & Peña, 1999). For  $u \geq 1$ , we have  $|1 \pm \gamma u| \leq |u \pm \gamma|$  since  $u^2(1 - \gamma^2) \geq 1 - \gamma^2$ . Then we have

$$\begin{aligned}
 &(p_0|1 + \gamma u|^q + p_0|1 - \gamma u|^q + (1 - 2p_0))^{\frac{1}{q}} \\
 &\leq (p_0|u|^q|1 + \gamma/u| + p_0|u|^q|1 - \gamma/u| + |u|^q(1 - 2p_0))^{\frac{1}{q}} \\
 &= |u| (p_0|1 + \gamma/u| + p_0|1 - \gamma/u| + (1 - 2p_0))^{\frac{1}{q}} \\
 &\leq |u| \left[ \frac{1}{2} |1 + \gamma/u|^q + |1 - \gamma/u|^q \right]^{\frac{1}{q}} \leq |u| \left[ \frac{1}{2} |1 + 1/u|^q + |1 - 1/u|^q \right]^{\frac{1}{q}} = (\mathbb{E}|1 + \epsilon u|^q)^{\frac{1}{q}}.
 \end{aligned}$$

where the last inequality is obtained by applying Eq. (3.2.4') in (Giné & Peña, 1999) and replacing  $u$  with  $1/u$ . Second, consider the case where  $x, y$  are vectors. Let  $z_1 = x + y$ ,  $z_2 = x - y$ ,  $u = \frac{\|z_1\| + \|z_2\|}{2}$  and  $v = \frac{\|z_1\| - \|z_2\|}{2}$ . In the second step, let  $\kappa = \frac{v}{u}$ , we have

$$\begin{aligned}
 (\mathbb{E}\|x + \gamma\sigma y\|^q)^{\frac{1}{q}} &= [p_0\|x + \gamma y\|^q + p_0\|x - \gamma y\|^q + (1 - 2p_0)\|x\|^q]^{\frac{1}{q}} \\
 &\leq \left[ p_0 \left( \frac{1 + \gamma}{2} \|z_1\| + \frac{1 - \gamma}{2} \|z_2\| \right)^q + p_0 \left( \frac{1 - \gamma}{2} \|z_1\| + \frac{1 + \gamma}{2} \|z_2\| \right)^q + (1 - 2p_0)\|x\|^q \right]^{\frac{1}{q}} \\
 &= [p_0 |u + \gamma v|^q + p_0 |u - \gamma v|^q + (1 - 2p_0)\|x\|^q]^{\frac{1}{q}} \\
 &\leq [p_0 |u + \gamma v|^q + p_0 |u - \gamma v|^q + (1 - 2p_0)u^q]^{\frac{1}{q}} \\
 &= |u| [p_0 |1 + \gamma\kappa|^q + p_0 |1 - \gamma\kappa|^q + (1 - 2p_0)]^{\frac{1}{q}} \\
 &\leq |u| \left[ \frac{|1 + \kappa|^p + |1 - \kappa|^p}{2} \right]^{\frac{1}{p}} = \left[ \frac{|u + v|^p + |u - v|^p}{2} \right]^{\frac{1}{p}},
 \end{aligned}$$

where we use the result for the the case that  $n = 1$  to obtain the last inequality. The second inequality is due to

$$\|x\| = \left\| \frac{1}{2}(x + y) + \frac{1}{2}(x - y) \right\| \leq \frac{1}{2}\|z_1\| + \frac{1}{2}\|z_2\| = u.$$

Third, we use induction to obtain the final result. Following (Giné & Peña, 1999), we only show  $n = 1$  implies  $n = 2$ .

Denote by  $\mu$  the measure on the probability space, by Fubini theorem:

$$\begin{aligned}
 & \left[ \mathbb{E}_{\sigma_1, \sigma_2} \left\| x + \gamma x_1 \sigma_1 + \gamma x_2 \sigma_2 + \gamma^2 x_{12} \sigma_1 \sigma_2 \right\|_2^q \right]^{\frac{1}{q}} \\
 &= \left[ \int \left( \int \left\| x + \gamma x_1 \sigma_1 + \gamma x_2 \sigma_2 + \gamma^2 x_{12} \sigma_1 \sigma_2 \right\|_2^q d\mu(\sigma_2) \right) d\mu(\sigma_1) \right]^{\frac{1}{q}} \\
 &\leq \left[ \int \left( \int \left\| x + \gamma x_1 \sigma_1 + x_2 \epsilon_2 + \gamma x_{12} \sigma_1 \epsilon_2 \right\|_2^p d\mu(\epsilon_2) \right)^{\frac{q}{p}} d\mu(\sigma_1) \right]^{\frac{1}{q}} \\
 &\leq \left[ \int \left( \int \left\| x + \gamma x_1 \sigma_1 + x_2 \epsilon_2 + \gamma x_{12} \sigma_1 \epsilon_2 \right\|_2^q d\mu(\sigma_1) \right)^{\frac{p}{q}} d\mu(\epsilon_2) \right]^{\frac{1}{p}} \\
 &\leq \left[ \int \int \left\| x + \gamma x_1 \sigma_1 + x_2 \epsilon_2 + x_{12} \sigma_1 \epsilon_2 \right\|_2^p d\mu(\epsilon_1) d\mu(\epsilon_2) \right]^{\frac{1}{p}} \\
 &= \left[ \mathbb{E}_{\epsilon_1, \epsilon_2} \left\| x + x_1 \sigma_1 + x_2 \sigma_2 + x_{12} \sigma_1 \sigma_2 \right\|_2^p \right]^{\frac{1}{p}},
 \end{aligned}$$

where the first and the third inequality is due to the induction hypothesis, and the second inequality is due to the Minkowski inequality. The remaining step is similar to the proof process in (Ying & Campbell, 2010) that bound the transductive Rademacher complexity by chaining technique. For  $j \in \mathbb{N}$ , let  $\alpha_k = 2^{-k}D$ . Denote by  $T_k$  the minimal  $\alpha_k$ -cover of  $\mathcal{F}_{\mathcal{W}}$  and  $f(\mathbf{w}^k; z, z')[\mathbf{w}]$  the element in  $T_k$  that covers  $f(\mathbf{w}; z, z')$ . Specifically, since  $\{f(\mathbf{w}^{(1)}; z, z')\}$  is a  $D$ -cover of  $\mathcal{W}$ , we set  $f(\mathbf{w}^0; z, z')[\mathbf{w}] = f(\mathbf{w}^{(1)}; z, z')$ . For arbitrary  $N \in \mathbb{N}$ :

$$\begin{aligned}
 & \mathbb{E}_{\sigma} \left[ \sup_{\mathbf{w} \in \mathcal{W}} \frac{1}{n} \sum_{1 \leq i < j \leq n} \sigma_i \sigma_j f(\mathbf{w}; z_i, z_j) \right] \\
 &= \mathbb{E}_{\sigma} \left[ \sup_{\mathbf{w} \in \mathcal{W}} \frac{1}{n} \left( \sum_{1 \leq i < j \leq n} \left( \sigma_i \sigma_j (f(\mathbf{w}; z_i, z_j) - f(\mathbf{w}^N; z_i, z_j))[\mathbf{w}] + \sum_{k=1}^N \sigma_i \sigma_j (f(\mathbf{w}^k; z_i, z_j)[\mathbf{w}] - f(\mathbf{w}^{k-1}; z_i, z_j)[\mathbf{w}]) \right. \right. \right. \\
 & \quad \left. \left. \left. + \sigma_i \sigma_j f(\mathbf{w}^{(1)}; z_i, z_j) \right) \right) \right] \\
 &\leq \mathbb{E}_{\sigma} \left[ \sup_{\mathbf{w} \in \mathcal{W}} \left( \frac{1}{n} \sum_{1 \leq i < j \leq n} \sigma_i \sigma_j (f(\mathbf{w}; z_i, z_j) - f(\mathbf{w}^N; z_i, z_j)[\mathbf{w}]) \right) \right] + \mathbb{E}_{\sigma} \left[ \frac{1}{n} \sum_{1 \leq i < j \leq n} \sigma_i \sigma_j f(\mathbf{w}^{(1)}; z_i, z_j) \right] \\
 & \quad + \sum_{k=1}^N \mathbb{E}_{\epsilon} \left[ \sup_{\mathbf{w} \in \mathcal{W}} \left( \frac{1}{n} \sum_{1 \leq i < j \leq n} \sigma_i \sigma_j (f(\mathbf{w}^k; z_i, z_j)[\mathbf{w}] - f(\mathbf{w}^{k-1}; z_i, z_j)[\mathbf{w}]) \right) \right].
 \end{aligned} \tag{15}$$

For the first term, we apply Cauchy-Schwarz inequality and obtain

$$\begin{aligned}
 & \mathbb{E}_{\sigma} \left[ \sup_{\mathbf{w} \in \mathcal{W}} \left( \frac{1}{n} \sum_{1 \leq i < j \leq n} \sigma_i \sigma_j (f(\mathbf{w}; z_i, z_j) - f(\mathbf{w}^N; z_i, z_j)[\mathbf{w}]) \right) \right] \\
 &\leq \left( \mathbb{E}_{\sigma} \left[ \sum_{1 \leq i < j \leq n} \sigma_i^2 \sigma_j^2 \right] \right)^{\frac{1}{2}} \left( \sup_{\mathbf{w} \in \mathcal{W}} \frac{1}{n^2} \sum_{1 \leq i < j \leq n} (f(\mathbf{w}; z_i, z_j) - f(\mathbf{w}^N; z_i, z_j)[\mathbf{w}])^2 \right)^{\frac{1}{2}} \leq n\alpha_N.
 \end{aligned} \tag{16}$$

For the second term, by and Jensen's inequality and Eq. (14),

$$\begin{aligned} \mathbb{E}_{\sigma} \left[ \frac{1}{n} \sum_{1 \leq i < j \leq n} \sigma_i \sigma_j f(\mathbf{w}^{(1)}; z_i, z_j) \right] &\leq \left( \mathbb{E}_{\epsilon} \left[ \left| \frac{1}{n} \sum_{1 \leq i < j \leq n} \epsilon_i \epsilon_j f(\mathbf{w}^{(1)}; z_i, z_j) \right|^2 \right] \right)^{\frac{1}{2}} \\ &= \left( \frac{1}{n^2} \sum_{1 \leq i < j \leq n} f^2(\mathbf{w}^{(1)}; z_i, z_j) \right)^{\frac{1}{2}} \leq M_0. \end{aligned} \quad (17)$$

Now we handle the last term. By Jensen's inequality,

$$\begin{aligned} &\exp \left\{ \lambda \mathbb{E}_{\sigma} \left[ \sup_{\mathbf{w} \in \mathcal{W}} \left| \frac{1}{n} \sum_{1 \leq i < j \leq n} \sigma_i \sigma_j (f(\mathbf{w}^k; z_i, z_j)[\mathbf{w}] - f(\mathbf{w}^{k-1}; z_i, z_j)[\mathbf{w}]) \right| \right] \right\} - 1 \\ &\leq \mathbb{E}_{\sigma} \left[ \exp \left\{ \lambda \sup_{\mathbf{w} \in \mathcal{W}} \left| \frac{1}{n} \sum_{1 \leq i < j \leq n} \sigma_i \sigma_j (f(\mathbf{w}^k; z_i, z_j)[\mathbf{w}] - f(\mathbf{w}^{k-1}; z_i, z_j)[\mathbf{w}]) \right| \right\} - 1 \right] \\ &= \mathbb{E}_{\sigma} \left[ \sup_{\mathbf{w} \in \mathcal{W}} \left( \exp \left\{ \lambda \left| \frac{1}{n} \sum_{1 \leq i < j \leq n} \sigma_i \sigma_j (f(\mathbf{w}^k; z_i, z_j)[\mathbf{w}] - f(\mathbf{w}^{k-1}; z_i, z_j)[\mathbf{w}]) \right| \right\} - 1 \right) \right] \\ &\leq |T_k| |T_{k-1}| \mathbb{E}_{\sigma} \left[ \left( \exp \left\{ \lambda \left| \frac{1}{n} \sum_{1 \leq i < j \leq n} \sigma_i \sigma_j (f(\mathbf{w}^k; z_i, z_j) - f(\mathbf{w}^{k-1}; z_i, z_j)) \right| \right\} - 1 \right) \right]. \end{aligned} \quad (18)$$

For any  $f(\mathbf{w}^{k-1}; z, z') \in T_{k-1}$  and  $f(\mathbf{w}^k; z, z') \in T_k$ , we have

$$\begin{aligned} &\mathbb{E}_{\sigma} \left[ \left( \exp \left\{ \lambda \left| \frac{1}{n} \sum_{1 \leq i < j \leq n} \sigma_i \sigma_j (f(\mathbf{w}^k; z_i, z_j) - f(\mathbf{w}^{k-1}; z_i, z_j)) \right| \right\} - 1 \right) \right] \\ &= \sum_{s \geq 1} \frac{1}{s!} \lambda^s \mathbb{E}_{\sigma} \left[ \left| \frac{1}{n} \sum_{1 \leq i < j \leq n} \sigma_i \sigma_j (f(\mathbf{w}^k; z_i, z_j) - f(\mathbf{w}^{k-1}; z_i, z_j)) \right|^s \right] \\ &\leq \sum_{s \geq 1} \frac{1}{s!} \lambda^s s^s \left[ \mathbb{E}_{\epsilon} \left| \frac{1}{n} \sum_{1 \leq i < j \leq n} \epsilon_i \epsilon_j (f(\mathbf{w}^k; z_i, z_j) - f(\mathbf{w}^{k-1}; z_i, z_j)) \right|^2 \right]^{\frac{s}{2}} \\ &\leq \sum_{s \geq 1} \left( e \lambda \left[ \mathbb{E}_{\epsilon} \left| \frac{1}{n} \sum_{1 \leq i < j \leq n} \epsilon_i \epsilon_j (f(\mathbf{w}^k; z_i, z_j) - f(\mathbf{w}^{k-1}; z_i, z_j)) \right|^2 \right]^{\frac{1}{2}} \right)^s, \end{aligned} \quad (19)$$

where the first inequality is by Eq. (14), and the second inequality is due to  $e^{-s} s^s \leq s!$ . Let

$$\lambda = \left( 2e \max_{f(\mathbf{w}^{k-1}; z, z') \in T_{k-1}, f(\mathbf{w}^k; z, z') \in T_k} \left[ \mathbb{E}_{\epsilon} \left| \frac{1}{n} \sum_{1 \leq i < j \leq n} \epsilon_i \epsilon_j (f(\mathbf{w}^k; z_i, z_j) - f(\mathbf{w}^{k-1}; z_i, z_j)) \right|^2 \right]^{\frac{1}{2}} \right)^{-1},$$

we obtain  $\mathbb{E}_{\epsilon} \left[ \left( \exp \left\{ \lambda \left| \frac{1}{n} \sum_{1 \leq i < j \leq n} \sigma_i \sigma_j (f(\mathbf{w}^k; z_i, z_j) - f(\mathbf{w}^{k-1}; z_i, z_j)) \right| \right\} - 1 \right) \right] \leq 1$ . Plugging this into Eq. (18)

yields

$$\begin{aligned}
 & \mathbb{E}_\epsilon \left[ \sup_{\mathbf{w} \in \mathcal{W}} \left| \frac{1}{n} \sum_{1 \leq i < j \leq n} \sigma_i \sigma_j (f(\mathbf{w}^k; z_i, z_j)[\mathbf{w}] - f(\mathbf{w}^{k-1}; z_i, z_j)[\mathbf{w}]) \right| \right] \\
 & \leq 2e \log(1 + |T_k| |T_{k-1}|) \max_{f(\mathbf{w}^{k-1}; z, z') \in T_{k-1}, f(\mathbf{w}^k; z, z') \in T_k} \left[ \mathbb{E}_\epsilon \left| \frac{1}{n} \sum_{1 \leq i < j \leq n} \epsilon_i \epsilon_j (f(\mathbf{w}^k; z_i, z_j) - f(\mathbf{w}^{k-1}; z_i, z_j)) \right|^2 \right]^{\frac{1}{2}} \\
 & \leq 2e \log(1 + |T_k| |T_{k-1}|) \max_{f(\mathbf{w}^{k-1}; z, z') \in T_{k-1}, f(\mathbf{w}^k; z, z') \in T_k} \left[ \frac{1}{n} \sum_{1 \leq i < j \leq n} [f(\mathbf{w}^k; z_i, z_j) - f(\mathbf{w}^{k-1}; z_i, z_j)]^2 \right]^{\frac{1}{2}} \\
 & = 2e \log(1 + |T_k| |T_{k-1}|) \left[ \frac{1}{n} \sum_{1 \leq i < j \leq n} [f(\mathbf{w}^k; z_i, z_j) - f(\mathbf{w}; z_i, z_j) + f(\mathbf{w}; z_i, z_j) - f(\mathbf{w}^{k-1}; z_i, z_j)]^2 \right]^{\frac{1}{2}} \\
 & \leq 2e \log(1 + |T_k| |T_{k-1}|) \left[ \frac{1}{n} \sum_{1 \leq i < j \leq n} [f(\mathbf{w}^k; z_i, z_j) - f(\mathbf{w}; z_i, z_j)]^2 \right]^{\frac{1}{2}} \\
 & \quad + 2e \log(1 + |T_k| |T_{k-1}|) \left[ \frac{1}{n} \sum_{1 \leq i < j \leq n} [f(\mathbf{w}; z_i, z_j) - f(\mathbf{w}^{k-1}; z_i, z_j)]^2 \right]^{\frac{1}{2}} \\
 & \leq 2e \log(1 + |T_k| |T_{k-1}|) [\alpha_{k-1} + \alpha_k] = 6e \log(1 + |T_k| |T_{k-1}|) \alpha_k.
 \end{aligned} \tag{20}$$

Plugging Eqs. (16, 17, 20) into Eq. (15), using the facts that  $\alpha_k = 2(\alpha_k - \alpha_{k+1})$  and  $|T_k| \geq |T_{k-1}|$ , we have

$$\begin{aligned}
 & \mathbb{E}_\sigma \left[ \sup_{\mathbf{w} \in \mathcal{W}} \frac{1}{n} \sum_{1 \leq i < j \leq n} \sigma_i \sigma_j f(\mathbf{w}; z_i, z_j) \right] \\
 & \leq M_0 + n\alpha_N + 6e \sum_{k=1}^N \log(1 + |T_k| |T_{k-1}|) \alpha_k = M_0 + n\alpha_N + 12e \sum_{k=1}^N \log(1 + |T_k| |T_{k-1}|) (\alpha_k - \alpha_{k+1}) \\
 & \leq M_0 + n\alpha_N + 24e \sum_{k=1}^N \log(1 + |T_k|) (\alpha_k - \alpha_{k+1}) \\
 & \leq M_0 + n\alpha_N + 24e \int_{\alpha_{N+1}}^{\alpha_0} \log(1 + \mathcal{N}(\alpha, \mathcal{F}_{\mathcal{W}}, d_{\mathcal{F}})) dr = M_0 + n\alpha_N + 24e \int_{\alpha_{N+1}}^D \log(1 + \mathcal{N}(\alpha, \mathcal{F}_{\mathcal{W}}, d_{\mathcal{F}})) dr.
 \end{aligned} \tag{21}$$

Taking the limit as  $N \rightarrow \infty$  we can obtain the result.

## B. Proof of Main Results

In this part, we present detailed proof of the results in the main body.

### B.1. Proof of Section 4.1

#### B.1.1. PROOF OF THEOREM 4.3

*Proof.* Following (El-Yaniv & Pechyony, 2007), let  $p = \frac{mu}{(m+u)^2}$ , the Transductive Rademacher complexity is defined as

$$\mathcal{R}_{m+u}(\mathbf{w}) = \left( \frac{1}{m} + \frac{1}{u} \right) \mathbb{E}_\sigma \left[ \sup_{\mathbf{w} \in \tilde{B}_R} \sum_{i=1}^{m+u} \sigma_i \ell(\mathbf{w}; z_i) \right],$$

where  $\sigma_i$  is a random variable taking value in  $\{\pm 1\}$  with probability  $p$  and 0 with probability  $1 - 2p$ . By Theorem 1 in (El-Yaniv & Pechyony, 2007), with probability at least  $1 - \delta/2$ ,

$$R_u(\mathbf{w}^{(T+1)}) \leq R_m(\mathbf{w}^{(T+1)}) + \mathcal{R}_{m+u}(\mathbf{w}) + c_0 Q \sqrt{\min(m, u)} + \sqrt{\frac{SQ}{2} \log \frac{2}{\delta}}, \quad (22)$$

where  $Q \triangleq (\frac{1}{m} + \frac{1}{u})$  and  $S \triangleq \frac{m+u}{(m+u-\frac{1}{2})(1-1/2(\max(m, u)))}$ .  $c_0 \triangleq \sqrt{\frac{32 \log(4e)}{3}}$  is a constant. Applying Lemma 1 in (El-Yaniv & Pechyony, 2007) with  $p_2 = \frac{1}{2}$ , we obtain

$$\mathcal{R}_{m+u}(\mathbf{w}) \leq \left(\frac{1}{m} + \frac{1}{u}\right) \mathbb{E}_\epsilon \left[ \sup_{\mathbf{w} \in B_R} \sum_{i=1}^{m+u} \epsilon_i \ell(\mathbf{w}; z_i) \right], \quad (23)$$

where  $\epsilon_i$  is the standard Rademacher random variable. Now we give an upper bound of the Transductive Rademacher Complexity by Dudley's integral technique. Denote by  $d_{\mathcal{H}_S}(\mathbf{w}, \tilde{\mathbf{w}}) = \left(\frac{1}{m+u} \sum_{i=1}^{m+u} [\ell(\mathbf{w}; z_i) - \ell(\tilde{\mathbf{w}}; z_i)]^2\right)^{\frac{1}{2}}$ . For  $j \in \mathbb{N}$ , let  $\alpha_j = 2^{-j}M$  with  $M = \sup_{\mathbf{w} \in B_R} d_{\mathcal{H}_S}(\mathbf{w}, \mathbf{w}^{(1)})$ . Denote by  $T_j$  the minimal  $\alpha_j$ -cover of  $B_R$  and  $\ell(\mathbf{w}^j; z)[\mathbf{w}]$  the element in  $T_j$  that covers  $\ell(\mathbf{w}; z)$ . Specifically, since  $\{\ell(\mathbf{w}^{(1)}; z)\}$  is a  $M$ -cover of  $B_R$ , we set  $\ell(\mathbf{w}^0; z)[\mathbf{w}] = \ell(\mathbf{w}^{(1)}; z)$  (recall that  $\mathbf{w}^{(1)}$  is the initialization parameter and  $\mathbf{w}^j$  is the associated parameter of  $\ell$  in  $T_j$ ). For arbitrary  $N \in \mathbb{N}$ :

$$\begin{aligned} & \mathbb{E}_\epsilon \left[ \sup_{\mathbf{w} \in B_R} \sum_{i=1}^{m+u} \epsilon_i \ell(\mathbf{w}; z_i) \right] \\ = & \mathbb{E}_\epsilon \left[ \sup_{\mathbf{w} \in B_R} \left( \sum_{i=1}^{m+u} \left( \epsilon_i (\ell(\mathbf{w}; z_i) - \ell(\mathbf{w}^N; z_i)[\mathbf{w}]) + \sum_{j=1}^N \epsilon_i (\ell(\mathbf{w}^j; z_i)[\mathbf{w}] - \ell(\mathbf{w}^{j-1}; z_i)[\mathbf{w}]) + \epsilon_i \ell(\mathbf{w}^{(1)}; z_i) \right) \right) \right] \\ \leq & \mathbb{E}_\epsilon \left[ \sup_{\mathbf{w} \in B_R} \left( \sum_{i=1}^{m+u} \epsilon_i (\ell(\mathbf{w}; z_i) - \ell(\mathbf{w}^N; z_i)[\mathbf{w}]) \right) \right] + \sum_{j=1}^N \mathbb{E}_\epsilon \left[ \sup_{\mathbf{w} \in B_R} \left( \sum_{i=1}^{m+u} \epsilon_i (\ell(\mathbf{w}^j; z_i)[\mathbf{w}] - \ell(\mathbf{w}^{j-1}; z_i)[\mathbf{w}]) \right) \right] \\ & + \mathbb{E}_\epsilon \left[ \sum_{i=1}^{m+u} \epsilon_i \ell(\mathbf{w}^{(1)}; z_i) \right]. \end{aligned} \quad (24)$$

For the first term, we apply Cauchy-Schwarz inequality and obtain

$$\begin{aligned} & \mathbb{E}_\epsilon \left[ \sup_{\mathbf{w} \in B_R} \left( \sum_{i=1}^{m+u} \epsilon_i (\ell(\mathbf{w}; z_i) - \ell(\mathbf{w}^N; z_i)[\mathbf{w}]) \right) \right] \\ \leq & \left( \mathbb{E}_\epsilon \left[ \sum_{i=1}^{m+u} \epsilon_i^2 \right] \right)^{\frac{1}{2}} \left( \sup_{\mathbf{w} \in B_R} \sum_{i=1}^{m+u} (\ell(\mathbf{w}; z_i) - \ell(\mathbf{w}^N; z_i)[\mathbf{w}])^2 \right)^{\frac{1}{2}} \leq (m+u) \alpha_N. \end{aligned} \quad (25)$$

By Massart's Lemma, we have

$$\mathbb{E}_\epsilon \left[ \sup_{\mathbf{w} \in B_R} \left( \sum_{i=1}^{m+u} \epsilon_i (\ell(\mathbf{w}^j; z_i)[\mathbf{w}] - \ell(\mathbf{w}^{j-1}; z_i)[\mathbf{w}]) \right) \right] \leq \sqrt{m+u} \sup_{\mathbf{w} \in B_R} d_{\mathcal{H}_S}(\mathbf{w}^j, \mathbf{w}^{j-1}) \sqrt{2 \log |T_j| |T_{j-1}|}. \quad (26)$$

By the Minkowski inequality,

$$\begin{aligned} & \sup_{\mathbf{w} \in B_R} d_{\mathcal{H}_S}(\mathbf{w}^j, \mathbf{w}^{j-1}) \\ = & \sup_{\mathbf{w} \in B_R} \left( \frac{1}{m+u} \sum_{i=1}^{m+u} [\ell(\mathbf{w}^j; z_i)[\mathbf{w}] - \ell(\mathbf{w}; z) + \ell(\mathbf{w}; z) - \ell(\mathbf{w}^{j-1}; z_i)[\mathbf{w}]]^2 \right)^{\frac{1}{2}} \\ \leq & \sup_{\mathbf{w} \in B_R} \left( \frac{1}{m+u} \sum_{i=1}^{m+u} [\ell(\mathbf{w}^j; z_i)[\mathbf{w}] - \ell(\mathbf{w}; z)]^2 \right)^{\frac{1}{2}} + \sup_{\mathbf{w} \in B_R} \left( \frac{1}{m+u} \sum_{i=1}^{m+u} [\ell(\mathbf{w}; z) - \ell(\mathbf{w}^{j-1}; z_i)[\mathbf{w}]]^2 \right)^{\frac{1}{2}} \\ = & \sup_{\mathbf{w} \in B_R} d_{\mathcal{H}_S}(\mathbf{w}^j, \mathbf{w}) + \sup_{\mathbf{w} \in B_R} d_{\mathcal{H}_S}(\mathbf{w}, \mathbf{w}^{j-1}) \leq \alpha_j + \alpha_{j-1} = 3\alpha_j. \end{aligned} \quad (27)$$

Plugging Eq. (27) into Eq. (26), using facts that  $\alpha_j = 2(\alpha_j - \alpha_{j+1})$  and  $|T_j| \geq |T_{j-1}|$ , taking summation over  $j$ ,

$$\begin{aligned}
 & \sum_{j=1}^N \mathbb{E}_\epsilon \left[ \sup_{\mathbf{w} \in B_R} \left( \sum_{i=1}^{m+u} \epsilon_i (\ell(\mathbf{w}^j; z_i)[\mathbf{w}] - \ell(\mathbf{w}^{j-1}; z_i)[\mathbf{w}]) \right) \right] \\
 & \leq 6\sqrt{m+u} \sum_{j=1}^N \alpha_j \sqrt{\log |T_j|} = 12\sqrt{m+u} \sum_{j=1}^N (\alpha_j - \alpha_{j+1}) \sqrt{\log |T_j|} \\
 & = 12\sqrt{m+u} \sum_{j=1}^N (\alpha_j - \alpha_{j+1}) \sqrt{\log \mathcal{N}(\alpha_j, \mathcal{H}_R, d_{\mathcal{H}_S})} \\
 & \leq 12\sqrt{m+u} \int_{\alpha_{N+1}}^{\alpha_0} \sqrt{\log \mathcal{N}(\alpha, \mathcal{H}_R, d_{\mathcal{H}_S})} d\alpha \leq 12\sqrt{m+u} \int_{\alpha_{N+1}}^{\infty} \sqrt{\log \mathcal{N}(\alpha, \mathcal{H}_R, d_{\mathcal{H}_S})} d\alpha.
 \end{aligned} \tag{28}$$

For the last term, by Khintchine-Kahane inequality (Latała & Oleszkiewicz, 1994),

$$\mathbb{E}_\epsilon \left[ \sum_{i=1}^{m+u} \epsilon_i \ell(\mathbf{w}^{(1)}; z_i) \right] \leq \left( \sum_{i=1}^{m+u} \ell^2(\mathbf{w}^{(1)}; z_i) \right)^{\frac{1}{2}} \leq b_\ell \sqrt{m+u}. \tag{29}$$

Taking the limit as  $N \rightarrow \infty$ , plugging Eq. (25), Eq. (28) and Eq. (29) into Eq. (24) and combining with Eq. (23) yield

$$\mathcal{R}_{m+u}(\mathbf{w}) \leq b_\ell \frac{(m+u)^{\frac{3}{2}}}{mu} + 12 \frac{(m+u)^{\frac{3}{2}}}{mu} \int_0^\infty \sqrt{\log \mathcal{N}(r, \mathcal{H}_R, d_{\mathcal{H}_S})} dr, \tag{30}$$

where  $\epsilon_i$  is the standard Rademacher random variable. Let  $\mathcal{H}_R = \{z \mapsto \ell(\mathbf{w}; z) \mid \mathbf{w} \in B_R\}$  be the parametric function space. One can verify that  $d_{\mathcal{H}_R}(\ell(\mathbf{w}; \cdot), \ell(\tilde{\mathbf{w}}; \cdot)) = \max_{z \in \mathcal{Z}} |\ell(\mathbf{w}; z) - \ell(\tilde{\mathbf{w}}; z)|$  is a metric in  $\mathcal{H}_R$ . we have

$$d_{\mathcal{H}_S} \leq \left( \frac{1}{m+u} \sum_{i=1}^{m+u} \left[ \max_{\mathbf{w}, \tilde{\mathbf{w}} \in B_R, z \in \mathcal{Z}} \ell(\mathbf{w}; z_i) - \ell(\tilde{\mathbf{w}}; z_i) \right]^2 \right)^{\frac{1}{2}} \leq d_{\mathcal{H}_R}.$$

By the definition of covering number, we have  $\mathcal{N}(r, \mathcal{H}_R, d_{\mathcal{H}_S}) \leq \mathcal{N}(r, \mathcal{H}_R, d_{\mathcal{H}_R})$ . Besides, applying Proposition 4.1 yields

$$d_{\mathcal{H}_R} = \max_{z \in \mathcal{Z}} |\ell(\mathbf{w}; z) - \ell(\tilde{\mathbf{w}}; z)| \leq L_{\mathcal{F}} \|\mathbf{w} - \tilde{\mathbf{w}}\|_2.$$

By the definition of covering number, we have  $\mathcal{N}(r, \mathcal{H}_R, d_{\mathcal{H}_R}) \leq \mathcal{N}\left(\frac{r}{L_{\mathcal{F}}}, B_R, d_2\right)$ . According to (Pisier, 1989),  $\log \mathcal{N}(r, B_R, d_2) \leq d \log(3R/r)$  holds. Therefore, we obtain

$$\log \mathcal{N}(r, \mathcal{H}_R, d_{\mathcal{H}_S}) \leq d \log \left( \frac{3L_{\mathcal{F}}R}{r} \right). \tag{31}$$

Furthermore,

$$d_{\mathcal{H}_S}^2(\mathbf{w}, \mathbf{w}^{(1)}) = \frac{1}{m+u} \sum_{i=1}^{m+u} \left[ \ell(\mathbf{w}; z_i) - \ell(\mathbf{w}^{(1)}; z_i) \right]^2 \leq L_{\mathcal{F}}^2 R^2,$$

where the last inequality is due to Proposition 4.1. This implies that

$$\int_0^\infty \sqrt{\log \mathcal{N}(r, \mathcal{H}_R, d_{\mathcal{H}_S})} dr = \int_0^{L_{\mathcal{F}}R} \sqrt{\log \mathcal{N}(r, \mathcal{H}_R, d_{\mathcal{H}_S})} dr. \tag{32}$$

Combining Eq. (30), Eq. (31), and Eq. (32) yields

$$\begin{aligned}
 \mathcal{R}_{m+u}(\mathbf{w}) & \leq 12 \frac{(m+u)^{\frac{3}{2}}}{mu} \sqrt{d} \int_0^{L_{\mathcal{F}}R} \sqrt{\log(3L_{\mathcal{F}}R/r)} dr \\
 & \leq 12 \frac{(m+u)^{\frac{3}{2}}}{mu} \sqrt{d} \left( \sqrt{\log 3} + \frac{3}{2} \sqrt{\pi} \right) L_{\mathcal{F}}R.
 \end{aligned} \tag{33}$$

Applying Theorem 47 in (Li & Liu, 2021) to bound  $R$  in Eq. (33) and plugging in Eq. (22) with probability  $1 - \delta/2$ , we conclude that with probability at least  $1 - \delta$ ,

$$R_u(\mathbf{w}^{(T+1)}) = \begin{cases} \mathcal{O}\left(L_{\mathcal{F}} \frac{(m+u)^{\frac{3}{2}}}{mu} \log^{\frac{1}{2}}(T) T^{\frac{1}{2}-\alpha} \log\left(\frac{1}{\delta}\right)\right) & \text{If } \alpha \in (0, \frac{1}{2}) \\ \mathcal{O}\left(L_{\mathcal{F}} \frac{(m+u)^{\frac{3}{2}}}{mu} \log(T) \log\left(\frac{1}{\delta}\right)\right) & \text{If } \alpha = \frac{1}{2} \\ \mathcal{O}\left(L_{\mathcal{F}} \frac{(m+u)^{\frac{3}{2}}}{mu} \log^{\frac{1}{2}}(T) \log\left(\frac{1}{\delta}\right)\right) & \text{If } \alpha \in (\frac{1}{2}, 1]. \end{cases}$$

### B.1.2. PROOF OF THEOREM 4.5

This proof extends the proof of Theorem 1 in literature (El-Yaniv & Pechyony, 2007) from scalar to vector. Let  $p = \frac{mu}{(m+u)^2}$ , we define the vector-valued Transductive Rademacher complexities:

$$\mathcal{R}_{m+u}(\mathbf{w}; p) = \mathbb{E}_{\sigma} \left[ \sup_{\mathbf{w} \in \mathcal{W}} \left\| \left( \frac{1}{m} + \frac{1}{u} \right) \sum_{i=1}^{m+u} \sigma_i \nabla \ell(\mathbf{w}; z_i) \right\|_2 \right],$$

where  $\sigma_i$  is a random variable taking value in  $\{\pm 1\}$  with probability  $p$  and 0 with probability  $1 - 2p$ . Following (El-Yaniv & Pechyony, 2007), we introduce the pairwise Rademacher variables  $\tilde{\sigma} = \{(\tilde{\sigma}_{i,1}, \tilde{\sigma}_{i,2})\}_{i=1}^{m+u}$  that satisfies:  $\mathbb{P}\{(\tilde{\sigma}_{i,1}, \tilde{\sigma}_{i,2}) = (\frac{1}{m}, \frac{1}{u})\} = \frac{mu}{(m+u)^2}$ ,  $\mathbb{P}\{(\tilde{\sigma}_{i,1}, \tilde{\sigma}_{i,2}) = (-\frac{1}{u}, -\frac{1}{m})\} = \frac{mu}{(m+u)^2}$ ,  $\mathbb{P}\{(\tilde{\sigma}_{i,1}, \tilde{\sigma}_{i,2}) = (\frac{1}{m}, -\frac{1}{m})\} = \frac{m^2}{(m+u)^2}$ ,  $\mathbb{P}\{(\tilde{\sigma}_{i,1}, \tilde{\sigma}_{i,2}) = (-\frac{1}{u}, \frac{1}{u})\} = \frac{u^2}{(m+u)^2}$ . It can be verify that

$$\mathcal{R}_{m+u}(\mathbf{w}; p) = \mathbb{E}_{\tilde{\sigma}} \left[ \sup_{\mathbf{w} \in \mathcal{W}} \left\| \sum_{i=1}^{m+u} (\sigma_{i,1} + \sigma_{i,2}) \nabla \ell(\mathbf{w}; z_i) \right\|_2 \right].$$

Denote by  $F_m(\mathbf{w}) \triangleq \frac{1}{m} \sum_{i=1}^m \nabla \ell(\mathbf{w}; z_i)$  and  $F_u(\mathbf{w}) \triangleq \frac{1}{u} \sum_{i=m+1}^{m+u} \nabla \ell(\mathbf{w}; z_i)$  the population gradient calculated on training samples and test samples. Let  $\Gamma : [m+u] \mapsto [m+u]$  be a symmetric group and  $\pi \in \Gamma$  a specific permutation on samples. We have

$$\begin{aligned} & \|F_m(\mathbf{w}; \pi) - F_u(\mathbf{w}; \pi)\|_2 \\ & \triangleq \sup_{\mathbf{w} \in \mathcal{W}} \left\| \frac{1}{m} \sum_{i=1}^m \nabla \ell(\mathbf{w}; z_{\pi(i)}) - \frac{1}{u} \sum_{i=m+1}^{m+u} \nabla \ell(\mathbf{w}; z_{\pi(i)}) \right\|_2 \\ & = \sup_{\mathbf{w} \in \mathcal{W}} \left\| \frac{1}{m} \sum_{i=1}^m \nabla \ell(\mathbf{w}; z_{\pi(i)}) - \frac{1}{m} \sum_{\pi' \in \Gamma} \sum_{i=1}^m \frac{\nabla \ell(\mathbf{w}; z_{\pi'(i)})}{(m+u)!} + \frac{1}{u} \sum_{\pi' \in \Gamma} \sum_{i=m+1}^{m+u} \frac{\nabla \ell(\mathbf{w}; z_{\pi'(i)})}{(m+u)!} - \frac{1}{u} \sum_{i=m+1}^{m+u} \nabla \ell(\mathbf{w}; z_{\pi(i)}) \right\|_2 \\ & = \sup_{\mathbf{w} \in \mathcal{W}} \left\| \sum_{\pi' \in \Gamma} \frac{1}{(m+u)!} \left[ \frac{1}{m} \sum_{i=1}^m (\nabla \ell(\mathbf{w}; z_{\pi(i)}) - \nabla \ell(\mathbf{w}; z_{\pi'(i)})) + \frac{1}{u} \sum_{i=m+1}^{m+u} (\nabla \ell(\mathbf{w}; z_{\pi'(i)}) - \nabla \ell(\mathbf{w}; z_{\pi(i)})) \right] \right\|_2 \\ & \leq \sum_{\pi' \in \Gamma} \frac{1}{(m+u)!} \sup_{\mathbf{w}} \left\| \frac{1}{m} \sum_{i=1}^m (\nabla \ell(\mathbf{w}; z_{\pi(i)}) - \nabla \ell(\mathbf{w}; z_{\pi'(i)})) + \frac{1}{u} \sum_{i=m+1}^{m+u} (\nabla \ell(\mathbf{w}; z_{\pi'(i)}) - \nabla \ell(\mathbf{w}; z_{\pi(i)})) \right\|_2 \\ & = \mathbb{E}_{\pi'} \sup_{\mathbf{w}} \left\| \frac{1}{m} \sum_{i=1}^m (\nabla \ell(\mathbf{w}; z_{\pi(i)}) - \nabla \ell(\mathbf{w}; z_{\pi'(i)})) + \frac{1}{u} \sum_{i=m+1}^{m+u} (\nabla \ell(\mathbf{w}; z_{\pi'(i)}) - \nabla \ell(\mathbf{w}; z_{\pi(i)})) \right\|_2 \\ & \triangleq \Phi(\pi). \end{aligned}$$

By Proposition 4.1,

$$\|\nabla \ell(\mathbf{w}; z)\|_2 \leq \left\| \nabla \ell(\mathbf{w}; z) - \nabla \ell(\mathbf{w}^{(1)}; z) \right\|_2 + \left\| \nabla \ell(\mathbf{w}^{(1)}; z) \right\|_2 \leq P_{\mathcal{F}} R + b_g. \quad (34)$$

where the second inequality is due to  $\|\mathbf{w} - \mathbf{w}^{(1)}\|_2 \leq R$  and  $R > 1$ . By Lemma 2 in (El-Yaniv & Pechyony, 2007), with probability at least  $1 - \delta/2$  over the random permutation  $\pi$ , for  $\mathbf{w} \in \mathcal{W}$ ,

$$\|F_m(\mathbf{w}) - F_u(\mathbf{w})\|_2 \leq \mathbb{E}_{\pi} [\Phi(\pi)] + (P_{\mathcal{F}} R + b_g) \sqrt{\frac{SQ}{2} \log \frac{2}{\delta}}. \quad (35)$$



Next we discuss how to give an upper bound of  $\mathbb{E}_\pi[\Phi(\pi)]$ . To achieve this, we have to build connection between  $\mathbb{E}_\pi[\Phi(\pi)]$  and  $\mathcal{R}_{m+u}(\mathbf{w}; p)$ . For a given permutation  $\pi$ , denote by  $\mathbf{a} \in \mathbb{R}^{m+u}$  a random vectors where  $a_i = \frac{1}{m}$  if  $i \in \{\pi(1), \dots, \pi(m)\}$  else  $-\frac{1}{u}$  and  $\mathbf{b} \in \mathbb{R}^{m+u}$  a random vectors where  $b_i = -\frac{1}{m}$  if  $i \in \{\pi(1), \dots, \pi(m)\}$  else  $\frac{1}{u}$ . To build the connection between  $\mathbf{a}$ ,  $\mathbf{b}$  and  $\tilde{\sigma}$ , an extra distribution that conditioned on  $\tilde{\sigma}$  is introduced. Denote by  $n_1(\tilde{\sigma}) \triangleq \sum_{i=1}^{m+u} \mathbb{I}\{(\tilde{\sigma}_{i,1}, \tilde{\sigma}_{i,2}) = (\frac{1}{m}, \frac{1}{u})\}$ ,  $n_2(\tilde{\sigma}) \triangleq \sum_{i=1}^{m+u} \mathbb{I}\{(\tilde{\sigma}_{i,1}, \tilde{\sigma}_{i,2}) = (\frac{1}{m}, -\frac{1}{m})\}$ , and  $n_3(\tilde{\sigma}) = \sum_{i=1}^{m+u} \mathbb{I}\{(\tilde{\sigma}_{i,1}, \tilde{\sigma}_{i,2}) = (-\frac{1}{u}, -\frac{1}{m})\}$  the random variables conditioned on  $\tilde{\sigma}$ , which indicate the number of pairs appearing in elements of  $\tilde{\sigma}$ . Let  $N_1(\tilde{\sigma}) = n_1(\tilde{\sigma}) + n_2(\tilde{\sigma})$  and  $N_2(\tilde{\sigma}) = n_2(\tilde{\sigma}) + n_3(\tilde{\sigma})$ , we denote by  $\mathfrak{R}(N_1, N_2)$  the distribution of  $\tilde{\sigma}$  conditioned on  $n_1, n_2$ , and  $n_3$ , which has fixed number of pairs and the randomness comes from permutations. Thus,  $\mathbf{a} + \mathbf{b}$  and  $\tilde{\sigma} \sim \mathfrak{R}(N_1(\tilde{\sigma}) = m, N_2(\tilde{\sigma}) = m)$  have the same distribution. Then we have

$$\begin{aligned} & \mathbb{E}_\pi[\Phi(\pi)] \\ &= \mathbb{E}_{\pi, \pi'} \sup_{\mathbf{w} \in \mathcal{W}} \left\| \frac{1}{m} \sum_{i=1}^m (\nabla \ell(\mathbf{w}; z_{\pi(i)}) - \nabla \ell(\mathbf{w}; z_{\pi'(i)})) + \frac{1}{u} \sum_{i=m+1}^{m+u} (\nabla \ell(\mathbf{w}; z_{\pi'(i)}) - \nabla \ell(\mathbf{w}; z_{\pi(i)})) \right\|_2 \\ &= \mathbb{E}_{\pi, \pi'} \sup_{\mathbf{w} \in \mathcal{W}} \left\| \sum_{i=1}^{m+u} (a_{\pi(i)} + b_{\pi'(i)}) \nabla \ell(\mathbf{w}; z_i) \right\|_2 \\ &= \mathbb{E}_{\tilde{\sigma} \sim \mathfrak{R}(m, m)} \sup_{\mathbf{w} \in \mathcal{W}} \left\| \sum_{i=1}^{m+u} (\tilde{\sigma}_{i,1} + \tilde{\sigma}_{i,2}) \nabla \ell(\mathbf{w}; z_i) \right\|_2. \end{aligned}$$

Denote by

$$\psi(N, N') = \mathbb{E}_{\tilde{\sigma} \sim \mathfrak{R}(N_1, N_2)} \left[ \sup_{\mathbf{w} \in \mathcal{W}} \left\| \sum_{i=1}^{m+u} (\sigma_{i,1} + \sigma_{i,2}) \nabla \ell(\mathbf{w}; z_i) \right\|_2 \right],$$

the Transductive Rademacher complexity where  $\tilde{\sigma}$  follows  $\mathfrak{R}(N, N')$  for given  $N_1$  and  $N_2$ . One can find that

$$\mathcal{R}_{m+u}(\mathbf{w}; p) = \mathbb{E}_{N_1(\tilde{\sigma}), N_2(\tilde{\sigma})} [\psi(N_1(\tilde{\sigma}), N_2(\tilde{\sigma}))]. \quad (36)$$

Besides, one can find that  $\mathbb{E}_{\tilde{\sigma}}[N_1(\tilde{\sigma})] = m$  and  $\mathbb{E}_{\tilde{\sigma}}[N_2(\tilde{\sigma})] = m$  hold. Therefore, we have

$$\mathbb{E}_\pi[\Phi(\pi)] = \psi(\mathbb{E}_{\tilde{\sigma}}[N_1(\tilde{\sigma})], \mathbb{E}_{\tilde{\sigma}}[N_2(\tilde{\sigma})]). \quad (37)$$

The last step is to give an upper bound of  $\psi(\mathbb{E}_{\tilde{\sigma}}[N_1(\tilde{\sigma})], \mathbb{E}_{\tilde{\sigma}}[N_2(\tilde{\sigma})]) - \mathbb{E}_{N_1(\tilde{\sigma}), N_2(\tilde{\sigma})} [\psi(N_1(\tilde{\sigma}), N_2(\tilde{\sigma}))]$ . Recall the definitions of  $\psi(N_1, N_2)$  and  $\psi(N'_1, N_2)$  are:

$$\begin{aligned} & \psi(N_1, N_2) \\ &= \mathbb{E}_{\pi, \pi'} \sup_{\mathbf{w} \in \mathcal{W}} \left\| \frac{1}{m} \sum_{i=1}^{N_1} \nabla \ell(\mathbf{w}; z_{\pi(i)}) - \frac{1}{m} \sum_{i=1}^{N_2} \nabla \ell(\mathbf{w}; z_{\pi'(i)}) + \frac{1}{u} \sum_{i=N_2+1}^{m+u} \nabla \ell(\mathbf{w}; z_{\pi'(i)}) - \frac{1}{u} \sum_{i=N_1+1}^{m+u} \nabla \ell(\mathbf{w}; z_{\pi(i)}) \right\|_2, \\ & \psi(N'_1, N_2) \\ &= \mathbb{E}_{\pi, \pi'} \sup_{\mathbf{w} \in \mathcal{W}} \left\| \frac{1}{m} \sum_{i=1}^{N'_1} \nabla \ell(\mathbf{w}; z_{\pi(i)}) - \frac{1}{m} \sum_{i=1}^{N_2} \nabla \ell(\mathbf{w}; z_{\pi'(i)}) + \frac{1}{u} \sum_{i=N_2+1}^{m+u} \nabla \ell(\mathbf{w}; z_{\pi'(i)}) - \frac{1}{u} \sum_{i=N'_1+1}^{m+u} \nabla \ell(\mathbf{w}; z_{\pi(i)}) \right\|_2. \end{aligned}$$

Without loss of generality, assume that  $N'_1 \leq N_1$ . Then we have

$$\begin{aligned} & |\psi(N_1, N_2) - \psi(N'_1, N_2)| \\ & \leq \mathbb{E}_\pi \sup_{\mathbf{w} \in \mathcal{W}} \left\| \left( \frac{1}{u} + \frac{1}{m} \right) \sum_{i=N'_1+1}^{N_1} \nabla \ell(\mathbf{w}; z_{\pi(i)}) \right\|_2 \leq |N_1 - N'_1| (PR + b_g) \left( \frac{1}{m} + \frac{1}{u} \right). \end{aligned} \quad (38)$$

Similarly, we have

$$\begin{aligned} & |\psi(N_1, N_2) - \psi(N_1, N'_2)| \\ & \leq \mathbb{E}_\pi \sup_{\mathbf{w} \in \mathcal{W}} \left\| \left( \frac{1}{u} + \frac{1}{m} \right) \sum_{i=N'_2+1}^{N_2} \nabla \ell(\mathbf{w}; z_{\pi(i)}) \right\|_2 \leq |N_2 - N'_2| (PR + b_g) \left( \frac{1}{m} + \frac{1}{u} \right). \end{aligned} \quad (39)$$

Combining Eq. (38), Eq. (39) and the inequality from (Devroye et al., 1996), we have

$$\begin{aligned}
 & \mathbb{P}_{N_1(\tilde{\sigma}), N_2(\tilde{\sigma})} \{ |\psi(N_1(\tilde{\sigma}), N_2(\tilde{\sigma})) - \psi(\mathbb{E}_{\tilde{\sigma}}[N_1(\tilde{\sigma})], \mathbb{E}_{\tilde{\sigma}}[N_2(\tilde{\sigma})])| \geq \epsilon \} \\
 & \leq \mathbb{P}_{N_1(\tilde{\sigma}), N_2(\tilde{\sigma})} \{ |\psi(N_1(\tilde{\sigma}), N_2(\tilde{\sigma})) - \psi(N_1, \mathbb{E}_{\tilde{\sigma}}[N_2(\tilde{\sigma})])| \geq \epsilon/2 \} \\
 & \quad + \mathbb{P}_{\tilde{\sigma}} \{ |\psi(N_1, \mathbb{E}_{\tilde{\sigma}}[N_2(\tilde{\sigma})]) - \psi(\mathbb{E}_{\tilde{\sigma}}[N_1(\tilde{\sigma})], \mathbb{E}_{\tilde{\sigma}}[N_2(\tilde{\sigma})])| \geq \epsilon/2 \} \\
 & \leq \mathbb{P}_{N_1(\tilde{\sigma}), N_2(\tilde{\sigma})} \{ |N_2(\tilde{\sigma}) - \mathbb{E}_{\tilde{\sigma}}[N_2(\tilde{\sigma})]|(PR + b_g)Q \geq \epsilon/2 \} + \mathbb{P}_{\tilde{\sigma}} \{ |N_1(\tilde{\sigma}) - \mathbb{E}_{\tilde{\sigma}}[N_1(\tilde{\sigma})]|(PR + b_g)Q \geq \epsilon/2 \} \\
 & \leq 4 \exp\{-3\epsilon^2/(32m(PR^\alpha + b_g)^2Q)\}.
 \end{aligned}$$

Applying the fact from Problem 12.1 in (Devroye et al., 1996), the following inequality holds

$$\mathbb{E}_{N_1(\tilde{\sigma}), N_2(\tilde{\sigma})} |\psi(N_1(\tilde{\sigma}), N_2(\tilde{\sigma})) - \psi(\mathbb{E}_{\tilde{\sigma}}[N_1(\tilde{\sigma})], \mathbb{E}_{\tilde{\sigma}}[N_2(\tilde{\sigma})])| \leq c_0(PR + b_g)Q\sqrt{\min(m, u)}, \quad (40)$$

where  $c_0 = \sqrt{\frac{32 \log(4e)}{3}}$ . Plugging Eq. (36), Eq. (37) and Eq. (40) into Eq. (35), with probability at least  $1 - \delta/2$ ,

$$\begin{aligned}
 & \|F_m(\mathbf{w}) - F_u(\mathbf{w})\|_2 \\
 & \leq \mathbb{E}_\pi[\Phi(\pi)] + (P_{\mathcal{F}}R + b_g)\sqrt{\frac{SQ}{2} \log \frac{2}{\delta}} \\
 & = \psi(\mathbb{E}_{\tilde{\sigma}}[N_1(\tilde{\sigma})], \mathbb{E}_{\tilde{\sigma}}[N_2(\tilde{\sigma})]) + (P_{\mathcal{F}}R + b_g)\sqrt{\frac{SQ}{2} \log \frac{2}{\delta}} \\
 & \leq \mathbb{E}_{N_1(\tilde{\sigma}), N_2(\tilde{\sigma})} [\psi(N_1(\tilde{\sigma}), N_2(\tilde{\sigma}))] + c_0(P_{\mathcal{F}}R + b_g)Q\sqrt{\min(m, u)} + (P_{\mathcal{F}}R + b_g)\sqrt{\frac{SQ}{2} \log \frac{2}{\delta}} \\
 & = \mathcal{R}_{m+u}(\mathbf{w}; p) + c_0(P_{\mathcal{F}}R + b_g)Q\sqrt{\min(m, u)} + (P_{\mathcal{F}}R + b_g)\sqrt{\frac{SQ}{2} \log \frac{2}{\delta}}.
 \end{aligned}$$

Till now, we have obtained the following inequality holds with probability at least  $1 - \delta/2$ :

$$\begin{aligned}
 & \sup_{\mathbf{w} \in B_R} \left\| \frac{1}{m} \sum_{i=1}^m \nabla \ell(\mathbf{w}; z_i) - \frac{1}{u} \sum_{i=m+1}^{m+u} \nabla \ell(\mathbf{w}; z_i) \right\|_2 \\
 & \leq \left( \frac{1}{m} + \frac{1}{u} \right) \mathbb{E}_\epsilon \left[ \sup_{\mathbf{w} \in B_R} \left\| \sum_{i=1}^{m+u} \epsilon_i \text{vec}(\nabla \ell(\mathbf{w}; z_i)) \right\|_2 \right] + c_0(PR + b_g)Q\sqrt{\min(m, u)} + (PR + b_g)\sqrt{\frac{S}{2} \left( \frac{1}{m} + \frac{1}{u} \right) \log \frac{1}{\delta}}. \quad (41)
 \end{aligned}$$

Let  $\mathcal{H}_R = \{(z, z') \mapsto \langle \nabla \ell(\mathbf{w}; z), \nabla \ell(\mathbf{w}; z') \rangle \mid \mathbf{w} \in B_R\}$  be the parametric function space, one can verify that

$$d_{\mathcal{H}_R} = \max_{z, z' \in \mathcal{Z}} |\langle \nabla \ell(\mathbf{w}; z), \nabla \ell(\mathbf{w}; z') \rangle - \langle \nabla \ell(\tilde{\mathbf{w}}; z), \nabla \ell(\tilde{\mathbf{w}}; z') \rangle|$$

is a metric in  $\mathcal{H}_R$ . Define

$$d_{\mathcal{H}_S}(\mathbf{w}, \tilde{\mathbf{w}}) = \left( \frac{1}{(m+u)^2} \sum_{1 \leq i < j \leq m+u} |\langle \nabla \ell(\mathbf{w}; z_i), \nabla \ell(\mathbf{w}; z_j) \rangle - \langle \nabla \ell(\tilde{\mathbf{w}}; z_i), \nabla \ell(\tilde{\mathbf{w}}; z_j) \rangle|^2 \right)^{\frac{1}{2}},$$

we have  $d_{\mathcal{H}_S}(\mathbf{w}, \tilde{\mathbf{w}}) \leq d_{\mathcal{H}_R}$ . By the definition of covering number, we have  $\mathcal{N}(r, \mathcal{H}_R, d_{\mathcal{H}_S}) \leq \mathcal{N}(r, \mathcal{H}_R, d_{\mathcal{H}_R})$ . Besides, applying Proposition 4.1 yields

$$\begin{aligned}
 & (m+u)^2 d_S^2(\mathbf{w}, \tilde{\mathbf{w}}) \\
 & = \sum_{1 \leq i < j \leq m+u} |\langle \nabla \ell(\mathbf{w}; z_i), \nabla \ell(\mathbf{w}; z_j) \rangle - \langle \nabla \ell(\tilde{\mathbf{w}}; z_i), \nabla \ell(\tilde{\mathbf{w}}; z_j) \rangle|^2 \\
 & \leq \sum_{1 \leq i < j \leq m+u} 2|\langle \nabla \ell(\mathbf{w}; z_i) - \nabla \ell(\tilde{\mathbf{w}}; z_i), \nabla \ell(\mathbf{w}; z_j) \rangle|^2 + 2|\langle \nabla \ell(\tilde{\mathbf{w}}; z_i), \nabla \ell(\mathbf{w}; z_j) - \nabla \ell(\tilde{\mathbf{w}}; z_j) \rangle|^2 \\
 & \leq \sum_{1 \leq i < j \leq m+u} 2\|\nabla \ell(\mathbf{w}; z_i) - \nabla \ell(\tilde{\mathbf{w}}; z_i)\|_2^2 \|\nabla \ell(\mathbf{w}; z_j)\|_2^2 + 2\|\nabla \ell(\tilde{\mathbf{w}}; z_i)\|_2^2 \|\nabla \ell(\mathbf{w}; z_j) - \nabla \ell(\tilde{\mathbf{w}}; z_j)\|_2^2 \\
 & \leq 2(m+u)(m+u-1)P_{\mathcal{F}}^2(P_{\mathcal{F}}R + b_g)^2 \max \left\{ \|\mathbf{w} - \tilde{\mathbf{w}}\|_2^{2\alpha}, \|\mathbf{w} - \tilde{\mathbf{w}}\|_2^2 \right\}.
 \end{aligned}$$

By the definition of covering number, we have

$$\mathcal{N}(r, \mathcal{H}_R, d_{\mathcal{H}_R}) \leq \mathcal{N} \left( \min \left\{ \left( \frac{r}{\sqrt{2}P_{\mathcal{F}}(P_{\mathcal{F}}R + b_g)} \right)^{\frac{1}{\alpha}}, \frac{r}{\sqrt{2}P_{\mathcal{F}}(P_{\mathcal{F}}R + b_g)} \right\}, B_R, d_2 \right).$$

According to (Pisier, 1989),  $\log \mathcal{N}(r, B_R, d_2) \leq d \log(3R/r)$  holds. Therefore, we obtain

$$\log \mathcal{N}(r, \mathcal{H}_R, d_{\mathcal{H}_S}) \leq \max \left\{ d \log \left( \frac{3R(\sqrt{2}P_{\mathcal{F}})^{\frac{1}{\alpha}} (P_{\mathcal{F}}R + b_g)^{\frac{1}{\alpha}}}{r^{\frac{1}{\alpha}}} \right), d \log \left( \frac{3\sqrt{2}P_{\mathcal{F}}R(P_{\mathcal{F}}R + b_g)}{r} \right) \right\}. \quad (42)$$

Denote by  $\frac{1}{m+u} \mathbb{E}_{\epsilon} \sum_{1 \leq i < j \leq m+u} \sigma_i \sigma_j h(z_i, z_j)$  the transductive Rademacher chaos complexity, we have

$$\begin{aligned} & \left( \mathbb{E}_{\sigma} \left[ \sup_{\mathbf{w} \in B_R} \left\| \sum_{i=1}^{m+u} \sigma_i \nabla \ell(\mathbf{w}; z_i) \right\|_2 \right] \right)^2 \\ & \leq \mathbb{E}_{\sigma} \left[ \sup_{\mathbf{w} \in B_R} \left\| \sum_{i=1}^{m+u} \sigma_i \nabla \ell(\mathbf{w}; z_i) \right\|_2^2 \right] \\ & = \mathbb{E}_{\sigma} \left[ \sup_{\mathbf{w} \in B_R} \sum_{i,j=1}^{m+u} \sigma_i \sigma_j \langle \nabla \ell(\mathbf{w}; z_i), \nabla \ell(\mathbf{w}; z_j) \rangle \right] \\ & = \mathbb{E}_{\sigma} \left[ \sup_{\mathbf{w} \in B_R} \sum_{i=1}^{m+u} \sigma_i^2 \|\nabla \ell(\mathbf{w}; z_i)\|_2^2 \right] + \mathbb{E}_{\sigma} \left[ \sup_{\mathbf{w} \in B_R} \sum_{i,j=1, i \neq j}^{m+u} \sigma_i \sigma_j \langle \nabla \ell(\mathbf{w}; z_i), \nabla \ell(\mathbf{w}; z_j) \rangle \right] \\ & \leq (m+u)(P_{\mathcal{F}}R + b_g)^2 + 2(m+u)\mathcal{U}(\mathcal{H}_R), \end{aligned} \quad (43)$$

Note that

$$\begin{aligned} & (m+u)^2 d_S^2(\mathbf{w}, \mathbf{w}^{(1)}) \\ & = \sum_{1 \leq i < j \leq m+u} |\langle \nabla \ell(\mathbf{w}; z_i), \nabla \ell(\mathbf{w}; z_j) \rangle - \langle \nabla \ell(\mathbf{w}^{(1)}; z_i), \nabla \ell(\mathbf{w}^{(1)}; z_j) \rangle| \\ & \leq \sum_{1 \leq i < j \leq m+u} 2 \left\| \nabla \ell(\mathbf{w}; z_i) - \nabla \ell(\mathbf{w}^{(1)}; z_i) \right\|_2^2 \|\nabla \ell(\mathbf{w}; z_j)\|_2^2 + 2 \left\| \nabla \ell(\mathbf{w}^{(1)}; z_i) \right\|_2^2 \left\| \nabla \ell(\mathbf{w}; z_j) - \nabla \ell(\mathbf{w}^{(1)}; z_j) \right\|_2^2 \\ & \leq 2(m+u)(m+u-1)P_{\mathcal{F}}^2 R^2 (P_{\mathcal{F}}R + b_g)^2. \end{aligned}$$

By Lemma A.7, we have  $\mathcal{U}(\mathcal{H}_R) \leq b_g^2 + 24e \int_0^{\sqrt{2}P_{\mathcal{F}}R(P_{\mathcal{F}}R + b_g)} \log(1 + \mathcal{N}(r, \mathcal{H}_R, d_{\mathcal{H}_S})) dr$ . Plugging in Eq. (42) yields

$$\begin{aligned} \mathcal{U}(\mathcal{H}_R) & \leq b_g^2 + 24e \int_0^{2P_{\mathcal{F}}R(P_{\mathcal{F}}R + b_g)} \log(1 + \mathcal{N}(r, \mathcal{H}_R, d_{\mathcal{H}_S})) dr \\ & \leq b_g^2 + 24e \int_0^{\sqrt{2}P_{\mathcal{F}}R(P_{\mathcal{F}}R + b_g)} (\log 2 + \log \mathcal{N}(r, \mathcal{H}_R, d_{\mathcal{H}_S})) dr \\ & = b_g^2 + 24\sqrt{2}eP_{\mathcal{F}}R(P_{\mathcal{F}}R + b_g) \log 2 + 24ed \int_0^{\sqrt{2}P_{\mathcal{F}}R(P_{\mathcal{F}}R + b_g)} \log \left( \frac{3R(2P_{\mathcal{F}})^{\frac{1}{\alpha}} (P_{\mathcal{F}}R + b_g)^{\frac{1}{\alpha}}}{r^{\frac{1}{\alpha}}} \right) dr \\ & \quad + 24ed \int_{\sqrt{2}P_{\mathcal{F}}R(P_{\mathcal{F}}R + b_g)}^{\sqrt{2}P_{\mathcal{F}}R(P_{\mathcal{F}}R + b_g)} \log \left( \frac{3\sqrt{2}P_{\mathcal{F}}R(P_{\mathcal{F}}R + b_g)}{r} \right) dr \\ & \leq b_g^2 + 24\sqrt{2}eP_{\mathcal{F}}R(P_{\mathcal{F}}R + b_g) \left[ d \log(3e^{\frac{1}{\alpha}} R) + dR \log(3e) + R \log 2 \right]. \end{aligned}$$

Applying Theorem 47 in (Li & Liu, 2021) to bound  $R$  in Eq. (33) with probability  $1 - \delta/2$  and combining it with Eqs. (43, 41), with probability at least  $1 - \delta$ ,

$$\begin{aligned} & \sup_{\mathbf{w} \in B_R} \left\| \frac{1}{m} \sum_{i=1}^m \nabla \ell(\mathbf{w}; z_i) - \frac{1}{u} \sum_{i=m+1}^{m+u} \nabla \ell(\mathbf{w}; z_i) \right\|_2 \\ &= \begin{cases} \mathcal{O}\left(\frac{(m+u)^{\frac{3}{2}}}{mu} \log^{\frac{1}{2}}(T) T^{\frac{1-2\alpha}{2}} \log\left(\frac{1}{\delta}\right)\right) & \text{if } \alpha \in (0, \frac{1}{2}) \\ \mathcal{O}\left(\frac{(m+u)^{\frac{3}{2}}}{mu} \log(T) \log\left(\frac{1}{\delta}\right)\right) & \text{if } \alpha = \frac{1}{2} \\ \mathcal{O}\left(\frac{(m+u)^{\frac{3}{2}}}{mu} \log^{\frac{1}{2}}(T) \log\left(\frac{1}{\delta}\right)\right) & \text{if } \alpha \in (\frac{1}{2}, 1]. \end{cases} \end{aligned}$$

### B.1.3. PROOF OF THEOREM 4.9

By Lemma 43 in (Li & Liu, 2021), we have

$$R_m(\mathbf{w}^{T+1}) - R_m(\hat{\mathbf{w}}^*) = \begin{cases} \mathcal{O}\left(\frac{1}{T^\alpha}\right) & \text{if } \alpha \in (0, 1) \\ \mathcal{O}\left(\frac{\log(T) \log^3(1/\delta)}{T}\right) & \text{if } \alpha = 1. \end{cases} \quad (44)$$

By Theorem 4.3,

$$R_u(\mathbf{w}^{(T+1)}) - R_m(\mathbf{w}^{(T+1)}) = \begin{cases} \mathcal{O}\left(L_{\mathcal{F}} \frac{(m+u)^{\frac{3}{2}}}{mu} \log^{\frac{1}{2}}(T) T^{\frac{1}{2}-\alpha} \log\left(\frac{1}{\delta}\right)\right) & \text{if } \alpha \in (0, \frac{1}{2}) \\ \mathcal{O}\left(L_{\mathcal{F}} \frac{(m+u)^{\frac{3}{2}}}{mu} \log(T) \log\left(\frac{1}{\delta}\right)\right) & \text{if } \alpha = \frac{1}{2} \\ \mathcal{O}\left(L_{\mathcal{F}} \frac{(m+u)^{\frac{3}{2}}}{mu} \log^{\frac{1}{2}}(T) \log\left(\frac{1}{\delta}\right)\right) & \text{if } \alpha \in (\frac{1}{2}, 1]. \end{cases} \quad (45)$$

Combing Eq. (44) and Eq. (45) yields the result.

## B.2. Proof of Section 4.2

### B.2.1. PROOF OF PROPOSITION 4.11

We first analyze the Lipschitz continuity. Denote by  $\mathbf{Z}^{(1)} = g(\tilde{\mathbf{A}})\mathbf{X}$ ,  $\mathbf{H}^{(1)} = \sigma(\mathbf{Z}^{(1)}\mathbf{W}_1)$  and  $\mathbf{Z}^{(2)} = g(\tilde{\mathbf{A}})\mathbf{H}^{(1)}$ , the forward process of GCN is given by  $\hat{\mathbf{Y}} = \text{Softmax}(\mathbf{Z}^{(2)}\mathbf{W}_2)$ . First, we have

$$\begin{aligned} \max_{i \in [n]} \|\mathbf{H}_{i^*}^{(1)}\|_2 &= \left\| \sigma \left( \sum_{j=1}^n [g(\tilde{\mathbf{A}})]_{ij} \mathbf{X}_{j^*} \mathbf{W}_1 \right) \right\|_2 \leq \left\| \sum_{j=1}^n [g(\tilde{\mathbf{A}})]_{ij} \mathbf{X}_{j^*} \mathbf{W}_1 \right\|_2 \\ &\leq \sum_{j=1}^n [g(\tilde{\mathbf{A}})]_{ij} \|\mathbf{X}_{j^*} \mathbf{W}_1\|_2 \leq c_X c_W \|g(\tilde{\mathbf{A}})\|_\infty, \end{aligned} \quad (46)$$

where the first inequality is due to the definition of  $\sigma(\cdot)$ . Similarly,

$$\max_{i \in [n]} \|\mathbf{Z}_{i^*}^{(1)}\|_2 = \left\| \sum_{j=1}^n [g(\tilde{\mathbf{A}})]_{ij} \mathbf{X}_{j^*} \right\|_2 \leq \sum_{j=1}^n [g(\tilde{\mathbf{A}})]_{ij} \|\mathbf{X}_{j^*}\|_2 \leq c_X \|g(\tilde{\mathbf{A}})\|_\infty \quad (47)$$

holds. Besides,

$$\max_{i \in [n]} \|\mathbf{Z}_{i^*}^{(2)}\|_2 = \left\| \sum_{j=1}^n [g(\tilde{\mathbf{A}})]_{ij} \mathbf{H}_{j^*}^{(1)} \right\|_2 \leq \sum_{j=1}^n [g(\tilde{\mathbf{A}})]_{ij} \|\mathbf{H}_{j^*}^{(1)}\|_2 \leq c_X c_W \|g(\tilde{\mathbf{A}})\|_\infty^2.$$

Then we analyze how  $\ell(\mathbf{W}_1, \mathbf{W}_2; z_i)$  change w.r.t.  $\mathbf{W}_2$  for fixed  $\mathbf{W}_1$  and  $i \in [n]$ :

$$\begin{aligned}
 & |\ell(\mathbf{W}_1, \mathbf{W}_2, z_i) - \ell(\mathbf{W}_1, \mathbf{W}'_2, z_i)| \\
 & \leq \sqrt{2} \left\| \mathbf{Z}_{i^*}^{(1)}(\mathbf{W}_2 - \mathbf{W}'_2) \right\|_2 = \sqrt{2} \left\| \sum_{j=1}^n [g(\tilde{\mathbf{A}})]_{ij} \mathbf{H}_{j^*}^{(1)}(\mathbf{W}_2 - \mathbf{W}'_2) \right\|_2 \\
 & \leq \sqrt{2} \sum_{j=1}^n [g(\tilde{\mathbf{A}})]_{ij} \left\| \mathbf{H}_{j^*}^{(1)} \right\|_2 \|\mathbf{W}_2 - \mathbf{W}'_2\| \leq \sqrt{2} \left\| g(\tilde{\mathbf{A}}) \right\|_\infty \max_{i \in [n]} \|\mathbf{H}_{i^*}^{(1)}\|_2 \|\mathbf{W}_2 - \mathbf{W}'_2\| \\
 & \leq c_X c_W \sqrt{2} \left\| g(\tilde{\mathbf{A}}) \right\|_\infty^2 \|\mathbf{W}_2 - \mathbf{W}'_2\| \leq c_X c_W \sqrt{2} \left\| g(\tilde{\mathbf{A}}) \right\|_\infty^2 \|\text{vec}[\mathbf{W}_2] - \text{vec}[\mathbf{W}'_2]\|_2,
 \end{aligned}$$

where the first inequality is due to the Lipschitz continuity property of softmax, and the last inequality is obtained by Eq. (46). Then we analyze the change of  $\ell(\mathbf{W}_1, \mathbf{W}_2; z_i)$  change w.r.t.  $\mathbf{W}_1$  for fixed  $\mathbf{W}_2$  and  $i \in [n]$ . Note that  $\mathbf{Z}^{(1)}$  and  $\mathbf{H}^{(1)}$  are function of  $\mathbf{W}_1$  in this case, which we denote by  $\mathbf{Z}^{(1)}(\mathbf{W}_1)$  and  $\mathbf{H}^{(1)}(\mathbf{W}_1)$ , respectively. Then,

$$\begin{aligned}
 & |\ell(\mathbf{W}_1, \mathbf{W}_2, z_i) - \ell(\mathbf{W}'_1, \mathbf{W}_2, z_i)| \\
 & \leq \sqrt{2} \left\| (\mathbf{Z}_{i^*}^{(1)}(\mathbf{W}_1) - \mathbf{Z}_{i^*}^{(1)}(\mathbf{W}'_1)) \mathbf{W}_2 \right\|_2 \\
 & \leq c_W \sqrt{2} \left\| \sum_{j=1}^n [g(\tilde{\mathbf{A}})]_{ij} (\mathbf{H}_{j^*}^{(1)}(\mathbf{W}_1) - \mathbf{H}_{j^*}^{(1)}(\mathbf{W}'_1)) \right\|_2 \\
 & \leq c_W \sqrt{2} \sum_{j=1}^n [g(\tilde{\mathbf{A}})]_{ij} \left\| \mathbf{Z}_{j^*}^{(1)}(\mathbf{W}_1 - \mathbf{W}'_1) \right\|_2 \\
 & \leq c_W \sqrt{2} \left\| g(\tilde{\mathbf{A}}) \right\|_\infty \max_{i \in [n]} \left\| \mathbf{Z}_{i^*}^{(1)} \right\|_2 \|\mathbf{W}_1 - \mathbf{W}'_1\| \leq c_X c_W \sqrt{2} \left\| g(\tilde{\mathbf{A}}) \right\|_\infty^2 \|\text{vec}[\mathbf{W}_1] - \text{vec}[\mathbf{W}'_1]\|_2.
 \end{aligned}$$

Let  $L_1 = L_2 = c_X c_W \sqrt{2} \left\| g(\tilde{\mathbf{A}}) \right\|_\infty^2$ , we conclude that  $|\ell(\mathbf{w}) - \ell(\mathbf{w}')| \leq L_{\mathcal{F}} \|\mathbf{w} - \mathbf{w}'\|_2$  holds with  $L_{\mathcal{F}} = 2c_X c_W \left\| g(\tilde{\mathbf{A}}) \right\|_\infty^2$ . By the chain rule, we have

$$\begin{aligned}
 \frac{\partial \ell(\mathbf{W}_1, \mathbf{W}_2; z_i)}{\partial \text{vec}[\mathbf{W}_2]} &= (\hat{\mathbf{y}}_i - \mathbf{y}_i) \otimes \mathbf{Z}_{i^*}^{(2)}, \\
 \frac{\partial \ell(\mathbf{W}_1, \mathbf{W}_2; z_i)}{\partial \text{vec}[\mathbf{W}_1]} &= \sum_{j=1}^n [g(\tilde{\mathbf{A}})]_{ij} \left( \sigma'(\mathbf{Z}_{j^*}^{(1)} \mathbf{W}_1) \odot (\hat{\mathbf{y}}_i - \mathbf{y}_i) \mathbf{W}_2^\top \right) \otimes \mathbf{Z}_{j^*}^{(1)}.
 \end{aligned}$$

We first analyze how  $\frac{\partial \ell(\mathbf{W}_1, \mathbf{W}_2; z_i)}{\partial \text{vec}[\mathbf{W}_2]}$  change w.r.t.  $\mathbf{W}_1$  and  $\mathbf{W}_2$ . Note that

$$\begin{aligned}
 & \left\| \frac{\partial \ell(\mathbf{W}_1, \mathbf{W}_2; z_i)}{\partial \text{vec}[\mathbf{W}_2]}(\mathbf{W}_1) - \frac{\partial \ell(\mathbf{W}_1, \mathbf{W}_2; z_i)}{\partial \text{vec}[\mathbf{W}_2]}(\mathbf{W}'_1) \right\|_2 \\
 &= \|\hat{\mathbf{y}}_i - \mathbf{y}_i\|_2 \left\| \mathbf{Z}_{i^*}^{(2)}(\mathbf{W}_1) - \mathbf{Z}_{i^*}^{(2)}(\mathbf{W}'_1) \right\|_2 + \|\hat{\mathbf{y}}_i(\mathbf{W}_1) - \hat{\mathbf{y}}_i(\mathbf{W}'_1)\|_2 \left\| \mathbf{Z}_{i^*}^{(2)} \right\|_2 \\
 &\leq \left( \sqrt{2} + 2c_X c_W^2 \left\| g(\tilde{\mathbf{A}}) \right\|_\infty^2 \right) \left\| \mathbf{Z}_{i^*}^{(2)}(\mathbf{W}_1) - \mathbf{Z}_{i^*}^{(2)}(\mathbf{W}'_1) \right\|_2 \\
 &= \left( \sqrt{2} + 2c_X c_W^2 \left\| g(\tilde{\mathbf{A}}) \right\|_\infty^2 \right) \left\| \sum_{j=1}^n [g(\tilde{\mathbf{A}})]_{ij} (\mathbf{H}_{j^*}^{(1)}(\mathbf{W}_1) - \mathbf{H}_{j^*}^{(1)}(\mathbf{W}'_1)) \right\|_2 \\
 &\leq \left( \sqrt{2} \left\| g(\tilde{\mathbf{A}}) \right\|_\infty + 2c_X c_W^2 \left\| g(\tilde{\mathbf{A}}) \right\|_\infty^3 \right) \|\mathbf{W}_1 - \mathbf{W}'_1\| \max_{i \in [n]} \left\| \mathbf{Z}_{i^*}^{(1)} \right\|_2 \\
 &\leq \left( \sqrt{2} c_X \left\| g(\tilde{\mathbf{A}}) \right\|_\infty^2 + 2c_X c_W^2 \left\| g(\tilde{\mathbf{A}}) \right\|_\infty^4 \right) \|\text{vec}[\mathbf{W}_1] - \text{vec}[\mathbf{W}'_1]\|.
 \end{aligned}$$

Besides,

$$\begin{aligned} & \left\| \frac{\partial \ell(\mathbf{W}_1, \mathbf{W}_2; z_i)}{\partial \text{vec}[\mathbf{W}_2]}(\mathbf{W}_2) - \frac{\partial \ell(\mathbf{W}_1, \mathbf{W}_2; z_i)}{\partial \text{vec}[\mathbf{W}_2]}(\mathbf{W}'_2) \right\|_2 \\ & \leq \|\hat{\mathbf{y}}_i(\mathbf{W}_2) - \hat{\mathbf{y}}_i(\mathbf{W}'_2)\|_2 \max_{i \in [n]} \|\mathbf{Z}_{i^*}^{(2)}\|_2 \leq 2 \|\mathbf{Z}^{(1)}(\mathbf{W}_2 - \mathbf{W}'_2)\|_2 \max_{i \in [n]} \|\mathbf{Z}_{i^*}^{(2)}\|_2 \\ & \leq 2 \|\mathbf{W}_2 - \mathbf{W}'_2\| \max_{i \in [n]} \|\mathbf{Z}_{i^*}^{(2)}\|_2 \leq 2c_X c_W \|g(\mathbf{A})\|_\infty^2 \|\text{vec}[\mathbf{W}_2] - \text{vec}[\mathbf{W}'_2]\|. \end{aligned}$$

Denote by  $P_{21} = \sqrt{2}c_X \|g(\tilde{\mathbf{A}})\|_\infty^2 + 2c_X c_W^2 \|g(\tilde{\mathbf{A}})\|_\infty^4$ ,  $P_{22} = 2c_X c_W \|g(\mathbf{A})\|_\infty^2$ ,  $\tilde{P}_{21} = \tilde{P}_{22} = 0$ , we obtain that  $\left\| \frac{\partial \ell(\mathbf{w}; z_i)}{\partial \text{vec}[\mathbf{W}_2]} - \frac{\partial \ell(\mathbf{w}'; z_i)}{\partial \text{vec}[\mathbf{W}_2]} \right\|_2 \leq \sum_{i=1}^2 P_{2i} \|\text{vec}[\mathbf{W}_i] - \text{vec}[\mathbf{W}'_i]\|_2 + \tilde{P}_{2i} \|\text{vec}[\mathbf{W}_i] - \text{vec}[\mathbf{W}'_i]\|_2^{\alpha_{2i}}$ . Then we analyze how  $\frac{\partial \ell(\mathbf{W}_1, \mathbf{W}_2; z_i)}{\partial \text{vec}[\mathbf{W}_1]}$  change w.r.t.  $\mathbf{W}_1$  and  $\mathbf{W}_2$ . Note that

$$\begin{aligned} & \left\| \frac{\partial \ell(\mathbf{W}_1, \mathbf{W}_2; z_i)}{\partial \text{vec}[\mathbf{W}_1]}(\mathbf{W}_1) - \frac{\partial \ell(\mathbf{W}_1, \mathbf{W}_2; z_i)}{\partial \text{vec}[\mathbf{W}_1]}(\mathbf{W}'_1) \right\|_2 \\ & = \left\| \sum_{j=1}^n [g(\tilde{\mathbf{A}})]_{ij} \left( (\sigma'(\mathbf{Z}_{j^*}^{(1)} \mathbf{W}_1) - \sigma'(\mathbf{Z}_{j^*}^{(1)} \mathbf{W}'_1)) \odot (\hat{\mathbf{y}}_i - \mathbf{y}_i) \mathbf{W}_2^\top \right) \otimes \mathbf{Z}_{j^*}^{(1)} \right\|_2 \\ & \quad + \left\| \sum_{j=1}^n [g(\tilde{\mathbf{A}})]_{ij} \left( \sigma'(\mathbf{Z}_{j^*}^{(1)} \mathbf{W}_1) \odot ((\hat{\mathbf{y}}_i(\mathbf{W}_1) - \hat{\mathbf{y}}_i(\mathbf{W}'_1))) \mathbf{W}_2^\top \right) \otimes \mathbf{Z}_{j^*}^{(1)} \right\|_2 \\ & \leq \|g(\tilde{\mathbf{A}})\|_\infty \|(\hat{\mathbf{y}}_i - \mathbf{y}_i) \mathbf{W}_2^\top\|_2 \max_{j \in [n]} \|\mathbf{Z}_{j^*}^{(1)}\|_2 \|\sigma'(\mathbf{Z}_{j^*}^{(1)} \mathbf{W}_1) - \sigma'(\mathbf{Z}_{j^*}^{(1)} \mathbf{W}'_1)\|_2 \\ & \quad + \|g(\tilde{\mathbf{A}})\|_\infty \max_{j \in [n]} \|\mathbf{Z}_{j^*}^{(1)}\|_2 \|\hat{\mathbf{y}}_i(\mathbf{W}_1) - \hat{\mathbf{y}}_i(\mathbf{W}'_1)\|_2 \\ & \leq c_X c_W P \sqrt{|\mathcal{Y}|} \|g(\tilde{\mathbf{A}})\|_\infty^2 \max_{j \in [n]} \|\mathbf{Z}_{j^*}^{(1)}(\mathbf{W}_1 - \mathbf{W}'_1)\|_2^{\tilde{\alpha}} + c_X c_W \|g(\tilde{\mathbf{A}})\|_\infty^2 \max_{j \in [n]} \|\mathbf{Z}_{i^*}^{(2)}(\mathbf{W}_1) - \mathbf{Z}_{i^*}^{(2)}(\mathbf{W}'_1)\|_2 \\ & \leq c_X c_W P \sqrt{|\mathcal{Y}|} \|g(\tilde{\mathbf{A}})\|_\infty^2 \|\mathbf{W}_1 - \mathbf{W}'_1\|_2^{\tilde{\alpha}} \max_{j \in [n]} \|\mathbf{Z}_{j^*}^{(1)}\|_2 + c_X c_W \|g(\tilde{\mathbf{A}})\|_\infty^2 \max_{j \in [n]} \|\mathbf{Z}_{i^*}^{(2)}(\mathbf{W}_1) - \mathbf{Z}_{i^*}^{(2)}(\mathbf{W}'_1)\|_2 \\ & \leq c_X^{1+\tilde{\alpha}} c_W P \sqrt{|\mathcal{Y}|} \|g(\tilde{\mathbf{A}})\|_\infty^{2+\tilde{\alpha}} \|\text{vec}[\mathbf{W}_1] - \text{vec}[\mathbf{W}'_1]\|_2^{\tilde{\alpha}} + c_X^2 c_W \|g(\tilde{\mathbf{A}})\|_\infty^4 \|\text{vec}[\mathbf{W}_1] - \text{vec}[\mathbf{W}'_1]\|_2, \end{aligned}$$

where we use the fact that the absolute value of each element of  $\sigma'(\mathbf{Z}_{j^*}^{(1)} \mathbf{W}_1)$  is less than 1. Similarly,

$$\begin{aligned} & \left\| \frac{\partial \ell(\mathbf{W}_1, \mathbf{W}_2; z_i)}{\partial \text{vec}[\mathbf{W}_1]}(\mathbf{W}_2) - \frac{\partial \ell(\mathbf{W}_1, \mathbf{W}_2; z_i)}{\partial \text{vec}[\mathbf{W}_1]}(\mathbf{W}'_2) \right\|_2 \\ & = \left\| \sum_{j=1}^n [g(\tilde{\mathbf{A}})]_{ij} \left( \sigma'(\mathbf{Z}_{j^*}^{(1)} \mathbf{W}_1) \odot (\hat{\mathbf{y}}_i - \mathbf{y}_i) (\mathbf{W}_2 - \mathbf{W}'_2)^\top \right) \otimes \mathbf{Z}_{j^*}^{(1)} \right\|_2 \\ & \quad + \left\| \sum_{j=1}^n [g(\tilde{\mathbf{A}})]_{ij} \left( \sigma'(\mathbf{Z}_{j^*}^{(1)} \mathbf{W}_1) \odot ((\hat{\mathbf{y}}_i(\mathbf{W}_2) - \hat{\mathbf{y}}_i(\mathbf{W}'_2))) \mathbf{W}_2^\top \right) \otimes \mathbf{Z}_{j^*}^{(1)} \right\|_2 \\ & \leq \|g(\tilde{\mathbf{A}})\|_\infty \|(\hat{\mathbf{y}}_i - \mathbf{y}_i) (\mathbf{W}_2 - \mathbf{W}'_2)^\top\|_2 \max_{j \in [n]} \|\mathbf{Z}_{j^*}^{(1)}\|_2 + \|g(\tilde{\mathbf{A}})\|_\infty \max_{j \in [n]} \|\mathbf{Z}_{j^*}^{(1)}\|_2 \|\hat{\mathbf{y}}_i(\mathbf{W}_2) - \hat{\mathbf{y}}_i(\mathbf{W}'_2)\|_2 \\ & \leq \left( \sqrt{2}c_X \|g(\tilde{\mathbf{A}})\|_\infty^2 + 2c_X^2 c_W \|g(\tilde{\mathbf{A}})\|_\infty^4 \right) \|\text{vec}[\mathbf{W}_2] - \text{vec}[\mathbf{W}'_2]\|_2. \end{aligned}$$

Denote by

$$\begin{aligned} P_{11} &= c_X^{1+\tilde{\alpha}} c_W P \sqrt{|\mathcal{Y}|} \|g(\tilde{\mathbf{A}})\|_\infty^{2+\tilde{\alpha}}, \tilde{P}_{11} = c_X^2 c_W \|g(\tilde{\mathbf{A}})\|_\infty^4, \\ P_{12} &= \sqrt{2}c_X \|g(\tilde{\mathbf{A}})\|_\infty^2 + 2c_X^2 c_W \|g(\tilde{\mathbf{A}})\|_\infty^4, \tilde{P}_{12} = 0, \end{aligned}$$

we obtain  $\left\| \frac{\partial \ell(\mathbf{w}; z_i)}{\partial \text{vec}[\mathbf{W}_1]} - \frac{\partial \ell(\mathbf{w}'; z_i)}{\partial \text{vec}[\mathbf{W}_1]} \right\|_2 \leq \sum_{i=1}^2 P_{1i} \|\text{vec}[\mathbf{W}_i] - \text{vec}[\mathbf{W}'_i]\|_2 + \tilde{P}_{1i} \|\text{vec}[\mathbf{W}_i] - \text{vec}[\mathbf{W}'_i]\|_2^{\tilde{\alpha}}$ . By Lemma A.4, we conclude that  $\|\nabla \ell(\mathbf{w}) - \nabla \ell(\mathbf{w}')\|_2 \leq P_{\mathcal{F}} \max\{\|\mathbf{w} - \mathbf{w}'\|_2, \|\mathbf{w} - \mathbf{w}'\|_2^{\tilde{\alpha}}\}$  where  $\mathbf{w} = [\text{vec}[\mathbf{W}_1]; \text{vec}[\mathbf{W}_2]]$ .

### B.2.2. PROOF OF PROPOSITION 4.12

We first analyze the Lipschitz continuity. Denote by

$$\begin{aligned} \mathbf{H}^{(0)} &= \sigma(\mathbf{X}\mathbf{W}_0), \\ \mathbf{H}^{(1)} &= \sigma\left(\left((1 - \alpha_1)g(\tilde{\mathbf{A}})\mathbf{H}^{(0)} + \alpha_1\mathbf{H}^{(0)}\right)\left((1 - \beta_1)\mathbf{I} + \beta_1\mathbf{W}_1\right)\right), \\ \mathbf{H}^{(2)} &= \sigma\left(\left((1 - \alpha_2)g(\tilde{\mathbf{A}})\mathbf{H}^{(1)} + \alpha_2\mathbf{H}^{(0)}\right)\left((1 - \beta_2)\mathbf{I} + \beta_2\mathbf{W}_2\right)\right), \end{aligned}$$

the forward process of GCNII is given by  $\hat{\mathbf{Y}} = \text{Softmax}(\mathbf{H}^{(2)}\mathbf{W}_3)$ . First, we have

$$\max_{i \in [n]} \left\| \mathbf{H}_{i^*}^{(0)} \right\|_2 = \max_{i \in [n]} \|\sigma(\mathbf{X}_{i^*}\mathbf{W}_0)\|_2 \leq c_X c_W.$$

Similarly, for  $\ell = 1$  and  $\ell = 2$ , denote by  $C_\ell = (1 - \beta_\ell) + \beta_\ell c_W$ , we have

$$\begin{aligned} & \max_{i \in [n]} \left\| \mathbf{H}_{i^*}^{(\ell)} \right\|_2 \\ &= \max_{i \in [n]} \left\| \sigma \left( \sum_{j=1}^n \left( (1 - \alpha_\ell) \left[ g(\tilde{\mathbf{A}}) \right]_{ij} \mathbf{H}_{j^*}^{(\ell-1)} + \alpha_\ell \mathbf{H}_{i^*}^{(0)} \right) \left( (1 - \beta_\ell) \mathbf{I} + \beta_\ell \mathbf{W}_\ell \right) \right) \right\|_2 \\ &\leq \max_{i \in [n]} \left\{ (1 - \alpha_\ell) \left\| \sum_{j=1}^n \left[ g(\tilde{\mathbf{A}}) \right]_{ij} \mathbf{H}_{j^*}^{(\ell-1)} \left( (1 - \beta_\ell) \mathbf{I} + \beta_\ell \mathbf{W}_\ell \right) \right\|_2 + \alpha_\ell \left\| \mathbf{H}_{i^*}^{(0)} \left( (1 - \beta_\ell) \mathbf{I} + \beta_\ell \mathbf{W}_\ell \right) \right\|_2 \right\} \\ &\leq \max_{i \in [n]} \left\{ (1 - \alpha_\ell) \sum_{j=1}^n \left[ g(\tilde{\mathbf{A}}) \right]_{ij} \left\| \mathbf{H}_{j^*}^{(\ell-1)} \left( (1 - \beta_\ell) \mathbf{I} + \beta_\ell \mathbf{W}_\ell \right) \right\|_2 + \alpha_\ell \left\| \mathbf{H}_{i^*}^{(0)} \left( (1 - \beta_\ell) \mathbf{I} + \beta_\ell \mathbf{W}_\ell \right) \right\|_2 \right\} \\ &\leq (1 - \alpha_\ell) \left\| g(\tilde{\mathbf{A}}) \right\|_\infty \left\| (1 - \beta_\ell) \mathbf{I} + \beta_\ell \mathbf{W}_\ell \right\|_2 \max_{i \in [n]} \left\| \mathbf{H}_{i^*}^{(\ell-1)} \right\|_2 + \alpha_\ell \left\| (1 - \beta_\ell) \mathbf{I} + \beta_\ell \mathbf{W}_\ell \right\|_2 \max_{i \in [n]} \left\| \mathbf{H}_{i^*}^{(0)} \right\|_2 \\ &\leq (1 - \alpha_\ell) C_\ell \left\| g(\tilde{\mathbf{A}}) \right\|_\infty \max_{i \in [n]} \left\| \mathbf{H}_{i^*}^{(\ell-1)} \right\|_2 + \alpha_\ell c_X c_W C_\ell. \end{aligned}$$

Let  $B_1 = \max_{i \in [n]} \left\| \mathbf{H}_{i^*}^{(1)} \right\|_2$  and  $B_2 = \max_{i \in [n]} \left\| \mathbf{H}_{i^*}^{(2)} \right\|_2$ , we have obtain  $B_1 = c_X c_W C_1 \left( (1 - \alpha_1) \left\| g(\tilde{\mathbf{A}}) \right\|_\infty + \alpha_1 \right)$  and  $B_2 = (1 - \alpha_2) C_2 \left\| g(\tilde{\mathbf{A}}) \right\|_\infty B_1 + \alpha_2 c_X c_W C_2$ . Next, we analyze the change of  $\mathbf{H}^{(1)}$  and  $\mathbf{H}^{(2)}$  w.r.t.  $\mathbf{W}_0$ ,  $\mathbf{W}_1$  and  $\mathbf{W}_2$ :

**Part A.** Note that

$$\begin{aligned} \Delta_{11} &\triangleq \left\| \mathbf{H}_{i^*}^{(1)}(\mathbf{W}_1) - \mathbf{H}_{i^*}^{(1)}(\mathbf{W}'_1) \right\|_2 \\ &\leq \beta_1 \left\| \left( (1 - \alpha_1) \sum_{j=1}^n \left[ g(\tilde{\mathbf{A}}) \right]_{ij} \mathbf{H}_{j^*}^{(0)} + \alpha_1 \mathbf{H}_{i^*}^{(0)} \right) (\mathbf{W}_1 - \mathbf{W}'_1) \right\|_2 \\ &\leq \beta_1 \left( (1 - \alpha_1) \sum_{j=1}^n \left[ g(\tilde{\mathbf{A}}) \right]_{ij} \left\| \mathbf{H}_{j^*}^{(0)} \right\|_2 + \alpha_1 \left\| \mathbf{H}_{i^*}^{(0)} \right\|_2 \right) \|\mathbf{W}_1 - \mathbf{W}'_1\| \\ &\leq \beta_1 \left( (1 - \alpha_1) \left\| g(\tilde{\mathbf{A}}) \right\|_\infty + \alpha_1 \right) c_X c_W \|\mathbf{W}_1 - \mathbf{W}'_1\| \\ &\leq \beta_1 \left( (1 - \alpha_1) \left\| g(\tilde{\mathbf{A}}) \right\|_\infty + \alpha_1 \right) c_X c_W \|\text{vec}[\mathbf{W}_1] - \text{vec}[\mathbf{W}'_1]\|_2 = \frac{\beta_1 B_1}{C_1} \|\text{vec}[\mathbf{W}_1] - \text{vec}[\mathbf{W}'_1]\|_2. \end{aligned}$$

Similarly, we have

$$\begin{aligned}
 \Delta_{10} &\triangleq \left\| \mathbf{H}_{i^*}^{(1)}(\mathbf{W}_0) - \mathbf{H}_{i^*}^{(1)}(\mathbf{W}'_0) \right\|_2 \\
 &\leq \left\| (1 - \alpha_1) \sum_{j=1}^n \left[ g(\tilde{\mathbf{A}}) \right]_{ij} \left( \mathbf{H}_{j^*}^{(0)}(\mathbf{W}_0) - \mathbf{H}_{j^*}^{(0)}(\mathbf{W}'_0) \right) + \alpha_1 \left( \mathbf{H}_{i^*}^{(0)}(\mathbf{W}_0) - \mathbf{H}_{i^*}^{(0)}(\mathbf{W}'_0) \right) \right\|_2 \left\| (1 - \beta_1)\mathbf{I} + \beta_1 \mathbf{W}_1 \right\| \\
 &\leq C_1 \left( (1 - \alpha_1) \sum_{j=1}^n \left[ g(\tilde{\mathbf{A}}) \right]_{ij} \left\| \mathbf{H}_{j^*}^{(0)}(\mathbf{W}_0) - \mathbf{H}_{j^*}^{(0)}(\mathbf{W}'_0) \right\|_2 + \alpha_1 \left\| \mathbf{H}_{i^*}^{(0)}(\mathbf{W}_0) - \mathbf{H}_{i^*}^{(0)}(\mathbf{W}'_0) \right\|_2 \right) \\
 &\leq C_1 \left( (1 - \alpha_1) \left\| g(\tilde{\mathbf{A}}) \right\|_\infty + \alpha_1 \right) c_X \|\mathbf{W}_0 - \mathbf{W}'_0\| = \frac{B_1}{c_W} \|\mathbf{W}_0 - \mathbf{W}'_0\|.
 \end{aligned}$$

**Part B.** Note that

$$\begin{aligned}
 \Delta_{22} &\triangleq \left\| \mathbf{H}_{i^*}^{(2)}(\mathbf{W}_2) - \mathbf{H}_{i^*}^{(2)}(\mathbf{W}'_2) \right\|_2 \\
 &\leq \beta_2 \left\| \left( (1 - \alpha_2) \sum_{j=1}^n \left[ g(\tilde{\mathbf{A}}) \right]_{ij} \mathbf{H}_{j^*}^{(1)} + \alpha_2 \mathbf{H}_{i^*}^{(0)} \right) (\mathbf{W}_2 - \mathbf{W}'_2) \right\|_2 \\
 &\leq \beta_2 \left( (1 - \alpha_2) \sum_{j=1}^n \left[ g(\tilde{\mathbf{A}}) \right]_{ij} \left\| \mathbf{H}_{j^*}^{(1)} \right\|_2 + \alpha_2 \left\| \mathbf{H}_{i^*}^{(0)} \right\|_2 \right) \|\mathbf{W}_2 - \mathbf{W}'_2\| \\
 &\leq \beta_2 \left( (1 - \alpha_2) \left\| g(\tilde{\mathbf{A}}) \right\|_\infty B_1 + \alpha_2 c_X c_W \right) \|\mathbf{W}_2 - \mathbf{W}'_2\| \\
 &\leq \beta_2 \left( (1 - \alpha_2) \left\| g(\tilde{\mathbf{A}}) \right\|_\infty B_1 + \alpha_2 c_X c_W \right) \|\text{vec}[\mathbf{W}_2] - \text{vec}[\mathbf{W}'_2]\|_2 = \frac{\beta_2 B_2}{C_2} \|\text{vec}[\mathbf{W}_1] - \text{vec}[\mathbf{W}'_1]\|_2.
 \end{aligned}$$

Similarly,

$$\begin{aligned}
 \Delta_{21} &\triangleq \left\| \mathbf{H}_{i^*}^{(2)}(\mathbf{W}_1) - \mathbf{H}_{i^*}^{(2)}(\mathbf{W}'_1) \right\|_2 \\
 &\leq (1 - \alpha_2) \sum_{j=1}^n \left[ g(\tilde{\mathbf{A}}) \right]_{ij} \left\| \mathbf{H}_{j^*}^{(1)}(\mathbf{W}_1) - \mathbf{H}_{j^*}^{(1)}(\mathbf{W}'_1) \right\|_2 \left\| (1 - \beta_2)\mathbf{I} + \beta_2 \mathbf{W}_2 \right\|_2 \\
 &\leq (1 - \alpha_2) C_2 \left\| g(\tilde{\mathbf{A}}) \right\|_\infty \max_{i \in [n]} \left\| \mathbf{H}_{i^*}^{(1)}(\mathbf{W}_1) - \mathbf{H}_{i^*}^{(1)}(\mathbf{W}'_1) \right\|_2 \leq (1 - \alpha_2) \beta_1 \frac{B_1 C_2}{C_1} \left\| g(\tilde{\mathbf{A}}) \right\|_\infty \|\text{vec}[\mathbf{W}_1] - \text{vec}[\mathbf{W}'_1]\|_2.
 \end{aligned}$$

Besides,

$$\begin{aligned}
 \Delta_{20} &\triangleq \left\| \mathbf{H}_{i^*}^{(2)}(\mathbf{W}_0) - \mathbf{H}_{i^*}^{(2)}(\mathbf{W}'_0) \right\|_2 \\
 &\leq (1 - \alpha_2) \sum_{j=1}^n \left[ g(\tilde{\mathbf{A}}) \right]_{ij} \left\| \mathbf{H}_{j^*}^{(1)}(\mathbf{W}_0) - \mathbf{H}_{j^*}^{(1)}(\mathbf{W}'_0) \right\|_2 \left\| (1 - \beta_2)\mathbf{I} + \beta_2 \mathbf{W}_2 \right\|_2 \\
 &\quad + \alpha_2 \left\| \mathbf{H}_{i^*}^{(0)}(\mathbf{W}_0) - \mathbf{H}_{i^*}^{(0)}(\mathbf{W}'_0) \right\|_2 \left\| (1 - \beta_2)\mathbf{I} + \beta_2 \mathbf{W}_2 \right\|_2 \\
 &\leq \left( (1 - \alpha_2) \frac{B_1 C_2}{c_W} \left\| g(\tilde{\mathbf{A}}) \right\|_\infty + \alpha_2 c_X C_2 \right) \|\text{vec}[\mathbf{W}_0] - \text{vec}[\mathbf{W}'_0]\|_2 = \frac{B_2}{c_W} \|\text{vec}[\mathbf{W}_0] - \text{vec}[\mathbf{W}'_0]\|_2.
 \end{aligned}$$

Now we are ready to analyze the Lipschitz continuity and Holder smoothness. Note that

$$\begin{aligned}
 &|\ell(\mathbf{W}_0, \mathbf{W}_1, \mathbf{W}_2, \mathbf{W}_3; z_i) - \ell(\mathbf{W}_0, \mathbf{W}_1, \mathbf{W}_2, \mathbf{W}'_3; z_i)| \\
 &\leq \sqrt{2} \left\| \mathbf{H}_{i^*}^{(2)}(\mathbf{W}_3 - \mathbf{W}'_3) \right\|_2 \leq \sqrt{2} \max_{i \in [n]} \left\| \mathbf{H}_{i^*}^{(2)} \right\|_2 \|\mathbf{W}_3 - \mathbf{W}'_3\|_2 \\
 &\leq \sqrt{2} B_2 \|\mathbf{W}_3 - \mathbf{W}'_3\|_2 \leq \sqrt{2} B_2 \|\text{vec}[\mathbf{W}_3] - \text{vec}[\mathbf{W}'_3]\|_2.
 \end{aligned}$$



Since  $\mathbf{H}^{(2)}$  is a variable related to  $\mathbf{W}_2$ , we have

$$\begin{aligned} & |\ell(\mathbf{W}_0, \mathbf{W}_1, \mathbf{W}_2, \mathbf{W}_3; z_i) - \ell(\mathbf{W}_0, \mathbf{W}_1, \mathbf{W}'_2, \mathbf{W}_3; z_i)| \\ & \leq \sqrt{2} \left\| (\mathbf{H}_{i*}^{(2)}(\mathbf{W}_2) - \mathbf{H}_{i*}^{(2)}(\mathbf{W}'_2)) \mathbf{W}_3 \right\|_2 \\ & \leq c_W \sqrt{2} \Delta_{22} = \frac{\beta_2 B_2}{C_2} c_W \sqrt{2} \|\text{vec}[\mathbf{W}_2] - \text{vec}[\mathbf{W}'_2]\|_2. \end{aligned}$$

Similarly, since  $\mathbf{H}^{(1)}$  and  $\mathbf{H}^{(2)}$  are variables related to  $\mathbf{W}_1$ ,

$$\begin{aligned} & |\ell(\mathbf{W}_0, \mathbf{W}_1, \mathbf{W}_2, \mathbf{W}_3; z_i) - \ell(\mathbf{W}_0, \mathbf{W}'_1, \mathbf{W}_2, \mathbf{W}_3; z_i)| \\ & \leq \sqrt{2} \left\| (\mathbf{H}_{i*}^{(2)}(\mathbf{W}_1) - \mathbf{H}_{i*}^{(2)}(\mathbf{W}'_1)) \mathbf{W}_3 \right\|_2 \\ & \leq c_W \sqrt{2} \Delta_{21} = (1 - \alpha_2) \beta_1 c_W \frac{B_1 C_2}{C_1} \left\| g(\tilde{\mathbf{A}}) \right\|_\infty \sqrt{2} \|\text{vec}[\mathbf{W}_1] - \text{vec}[\mathbf{W}'_1]\|_2. \end{aligned}$$

Lastly, since  $\mathbf{H}^{(0)}$ ,  $\mathbf{H}^{(1)}$  and  $\mathbf{H}^{(2)}$  are variables related to  $\mathbf{W}_0$ ,

$$\begin{aligned} & |\ell(\mathbf{W}_0, \mathbf{W}_1, \mathbf{W}_2, \mathbf{W}_3; z_i) - \ell(\mathbf{W}'_0, \mathbf{W}_1, \mathbf{W}_2, \mathbf{W}_3; z_i)| \\ & \leq \sqrt{2} \left\| (\mathbf{H}_{i*}^{(2)}(\mathbf{W}_0) - \mathbf{H}_{i*}^{(2)}(\mathbf{W}'_0)) \mathbf{W}_3 \right\|_2 \\ & \leq c_W \sqrt{2} \Delta_{20} = \sqrt{2} B_2 \|\text{vec}[\mathbf{W}_0] - \text{vec}[\mathbf{W}'_0]\|_2. \end{aligned}$$

Denote by

$$L_{\mathcal{F}} = \sqrt{4B_2^2 + 2c_W^2 \frac{\beta_2^2 B_2^2}{C_2^2} + 2(1 - \alpha_2)^2 \beta_1^2 c_W^2 \frac{B_1^2 C_2^2}{C_1^2} \left\| g(\tilde{\mathbf{A}}) \right\|_\infty^2},$$

by Lemma A.4, we conclude that  $|\ell(\mathbf{w}) - \ell(\mathbf{w}')| \leq L_{\mathcal{F}} \|\mathbf{w} - \mathbf{w}'\|_2$  holds. Then we discuss the smoothness. By the chain rule, we have

$$\begin{aligned} & \frac{\partial \ell(\mathbf{W}_0, \mathbf{W}_1, \mathbf{W}_2, \mathbf{W}_3; z_i)}{\partial \text{vec}[\mathbf{W}_3]} = (\hat{\mathbf{y}}_i - \mathbf{y}_i) \otimes \mathbf{H}_{i*}^{(2)}, \\ & \frac{\partial \ell(\mathbf{W}_0, \mathbf{W}_1, \mathbf{W}_2, \mathbf{W}_3; z_i)}{\partial \text{vec}[\mathbf{W}_2]} = \alpha_2 \beta_2 \delta_i \otimes \mathbf{H}_{i*}^{(0)} + (1 - \alpha_2) \beta_2 \sum_{j=1}^n \left[ g(\tilde{\mathbf{A}}) \right]_{ij} \delta_i \otimes \mathbf{H}_{j*}^{(1)}, \\ & \frac{\partial \ell(\mathbf{W}_0, \mathbf{W}_1, \mathbf{W}_2, \mathbf{W}_3; z_i)}{\partial \text{vec}[\mathbf{W}_1]} = \alpha_1 \beta_1 \sum_{j=1}^n \left[ g(\tilde{\mathbf{A}}) \right]_{ij} \delta_{ij} \otimes \mathbf{H}_{j*}^{(0)} + (1 - \alpha_1) \beta_1 \sum_{j=1}^n \sum_{k=1}^n \left[ g(\tilde{\mathbf{A}}) \right]_{ij} \left[ g(\tilde{\mathbf{A}}) \right]_{jk} \delta_{ij} \otimes \mathbf{H}_{k*}^{(0)}, \\ & \frac{\partial \ell(\mathbf{W}_0, \mathbf{W}_1, \mathbf{W}_2, \mathbf{W}_3; z_i)}{\partial \text{vec}[\mathbf{W}_0]} \\ & = \alpha_2 ((\delta_i((1 - \beta_2)\mathbf{I} + \beta_2 \mathbf{W}_2^\top)) \odot \mathbf{H}_{i*}^{(0)}) \otimes \mathbf{X}_{i*} \\ & + \alpha_1 \sum_{j=1}^n \left[ g(\tilde{\mathbf{A}}) \right]_{ij} ((\delta_{ij}((1 - \beta_1)\mathbf{I} + \beta_1 \mathbf{W}_1^\top)) \odot \mathbf{H}_{j*}^{(0)}) \otimes \mathbf{X}_{j*} \\ & + (1 - \alpha_1) \sum_{j=1}^n \sum_{k=1}^n \left[ g(\tilde{\mathbf{A}}) \right]_{ij} \left[ g(\tilde{\mathbf{A}}) \right]_{jk} ((\delta_{ij}((1 - \beta_1)\mathbf{I} + \beta_1 \mathbf{W}_1^\top)) \odot \mathbf{H}_{k*}^{(0)}) \otimes \mathbf{X}_{k*}, \end{aligned}$$

where

$$\begin{aligned} \delta_i & = ((\hat{\mathbf{y}}_i - \mathbf{y}_i) \mathbf{W}_3^\top) \odot \sigma'(\mathbf{H}_{i*}^{(2)}), \\ \delta_{ij} & = (1 - \alpha_2) \sigma'(\mathbf{H}_{j*}^{(1)}) \odot (\delta_i((1 - \beta_2) \mathbf{W}_2^\top + \beta_2 \mathbf{I})). \end{aligned}$$

We first analyze how  $\delta_i$  and  $\delta_{ij}$  change w.r.t.  $\mathbf{W}_0$ ,  $\mathbf{W}_1$ ,  $\mathbf{W}_2$  and  $\mathbf{W}_3$ .

**Part C.** For  $i \in [n]$ , we have

$$\begin{aligned}
 & \|\delta_i(\mathbf{W}_3) - \delta_i(\mathbf{W}'_3)\|_2 \\
 & \leq \left\| (\hat{\mathbf{y}}_i - \mathbf{y}_i)(\mathbf{W}_3 - \mathbf{W}'_3)^\top \odot \sigma'(\mathbf{H}_{i*}^{(2)}) \right\|_2 + \left\| (\hat{\mathbf{y}}_i(\mathbf{W}_3) - \hat{\mathbf{y}}_i(\mathbf{W}'_3))\mathbf{W}'_3{}^\top \odot \sigma'(\mathbf{H}_{i*}^{(2)}) \right\|_2 \\
 & \leq \|(\hat{\mathbf{y}}_i - \mathbf{y}_i)(\mathbf{W}_3 - \mathbf{W}'_3)\|_2 + c_W \|(\hat{\mathbf{y}}_i(\mathbf{W}_3) - \hat{\mathbf{y}}_i(\mathbf{W}'_3))\|_2 \\
 & \leq \sqrt{2} \|\mathbf{W}_3 - \mathbf{W}'_3\|_2 + 2c_W \|\mathbf{W}_3 - \mathbf{W}'_3\|_2 \max_{i \in [n]} \|\mathbf{H}_{i*}^{(2)}\|_2 = (\sqrt{2} + 2c_W B_2) \|\text{vec}[\mathbf{W}_3] - \text{vec}[\mathbf{W}'_3]\|_2,
 \end{aligned} \tag{48}$$

where we use the fact that the absolute value of each component in  $\sigma'(\mathbf{H}_{i*}^{(2)})$  is less than 1. Similarly, we have

$$\begin{aligned}
 & \|\delta_i(\mathbf{W}_2) - \delta_i(\mathbf{W}'_2)\|_2 \\
 & \leq \left\| (\hat{\mathbf{y}}_i - \mathbf{y}_i)\mathbf{W}_3{}^\top \odot (\sigma'(\mathbf{H}_{i*}^{(2)})(\mathbf{W}_2) - \sigma'(\mathbf{H}_{i*}^{(2)})(\mathbf{W}'_2)) \right\|_2 + \left\| ((\hat{\mathbf{y}}_i(\mathbf{W}_2) - \hat{\mathbf{y}}_i(\mathbf{W}'_2))\mathbf{W}_3{}^\top) \odot \sigma'(\mathbf{H}_{i*}^{(2)}) \right\|_2 \\
 & \leq P \|(\hat{\mathbf{y}}_i - \mathbf{y}_i)\mathbf{W}_3{}^\top\|_2 \left\| \mathbf{H}_{i*}^{(2)}(\mathbf{W}_2) - \mathbf{H}_{i*}^{(2)}(\mathbf{W}'_2) \right\|_2^{\tilde{\alpha}} + 2c_W^2 \left\| \mathbf{H}_{i*}^{(2)}(\mathbf{W}_2) - \mathbf{H}_{i*}^{(2)}(\mathbf{W}'_2) \right\|_2 \\
 & \leq P\sqrt{2}c_W \Delta_{22}^{\tilde{\alpha}} + 2c_W^2 \Delta_{22} \\
 & = \sqrt{2}c_W P \left( \frac{\beta_2 B_2}{C_2} \right)^{\tilde{\alpha}} \|\text{vec}[\mathbf{W}_2] - \text{vec}[\mathbf{W}'_2]\|_2^{\tilde{\alpha}} + 2c_W^2 \frac{\beta_2 B_2}{C_2} \|\text{vec}[\mathbf{W}_2] - \text{vec}[\mathbf{W}'_2]\|_2.
 \end{aligned} \tag{49}$$

Besides,

$$\begin{aligned}
 & \|\delta_i(\mathbf{W}_1) - \delta_i(\mathbf{W}'_1)\|_2 \leq \sqrt{2}c_W P \Delta_{21}^{\tilde{\alpha}} + 2c_W^2 \Delta_{21} \\
 & = \sqrt{2}c_W P \left( (1 - \alpha_2)\beta_1 \frac{B_1 C_2}{C_1} \left\| g(\tilde{\mathbf{A}}) \right\|_\infty \right)^{\tilde{\alpha}} \|\text{vec}[\mathbf{W}_1] - \text{vec}[\mathbf{W}'_1]\|_2^{\tilde{\alpha}} \\
 & \quad + 2(1 - \alpha_2)\beta_1 c_W^2 \frac{B_1 C_2}{C_1} \left\| g(\tilde{\mathbf{A}}) \right\|_\infty \|\text{vec}[\mathbf{W}_1] - \text{vec}[\mathbf{W}'_1]\|_2, \\
 & \|\delta_i(\mathbf{W}_0) - \delta_i(\mathbf{W}'_0)\|_2 \leq \sqrt{2}c_W P \Delta_{20}^{\tilde{\alpha}} + 2c_W^2 \Delta_{20} \\
 & = \sqrt{2}c_W P \left( \frac{B_2}{c_W} \right)^{\tilde{\alpha}} \|\text{vec}[\mathbf{W}_0] - \text{vec}[\mathbf{W}'_0]\|_2^{\tilde{\alpha}} + 2c_W B_2 \|\text{vec}[\mathbf{W}_0] - \text{vec}[\mathbf{W}'_0]\|_2.
 \end{aligned} \tag{50}$$

**Part D.** For  $i \in [n]$ , we have

$$\begin{aligned}
 & \|\delta_{ij}(\mathbf{W}_3) - \delta_{ij}(\mathbf{W}'_3)\|_2 \\
 & = (1 - \alpha_2) \left\| \sigma'(\mathbf{H}_{j*}^{(1)}) \odot ((\delta_i(\mathbf{W}_3) - \delta_i(\mathbf{W}'_3))((1 - \beta_2)\mathbf{W}_2{}^\top + \beta_2\mathbf{I})) \right\|_2 \\
 & \leq (1 - \alpha_2) \left\| (\delta_i(\mathbf{W}_3) - \delta_i(\mathbf{W}'_3))((1 - \beta_2)\mathbf{W}_2{}^\top + \beta_2\mathbf{I}) \right\|_2 \\
 & \leq (1 - \alpha_2) C_2 \|\delta_i(\mathbf{W}_3) - \delta_i(\mathbf{W}'_3)\|_2 \leq (1 - \alpha_2) C_2 (\sqrt{2} + 2c_W B_2) \|\text{vec}[\mathbf{W}_3] - \text{vec}[\mathbf{W}'_3]\|_2.
 \end{aligned} \tag{51}$$

Similarly, we have

$$\begin{aligned}
 & \|\delta_{ij}(\mathbf{W}_2) - \delta_{ij}(\mathbf{W}'_2)\|_2 \\
 & \leq (1 - \alpha_2) \left\| \sigma'(\mathbf{H}_{j*}^{(1)}) \odot ((\delta_i(\mathbf{W}_2) - \delta_i(\mathbf{W}'_2))((1 - \beta_2)\mathbf{W}_2{}^\top + \beta_2\mathbf{I})) \right\|_2 \\
 & \quad + (1 - \alpha_2)(1 - \beta_2) \left\| \sigma'(\mathbf{H}_{j*}^{(1)}) \odot ((\delta_i(\mathbf{W}_2 - \mathbf{W}'_2)^\top)) \right\|_2 \\
 & \leq (1 - \alpha_2) \left\| (\delta_i(\mathbf{W}_2) - \delta_i(\mathbf{W}'_2))((1 - \beta_2)\mathbf{W}_2{}^\top + \beta_2\mathbf{I}) \right\|_2 + (1 - \alpha_2)(1 - \beta_2) \|\delta_i(\mathbf{W}_2 - \mathbf{W}'_2)^\top\|_2 \\
 & \leq (1 - \alpha_2) C_2 \|\delta_i(\mathbf{W}_2) - \delta_i(\mathbf{W}'_2)\|_2 + (1 - \alpha_2)(1 - \beta_2) \|\delta_i\|_2 \|\mathbf{W}_2 - \mathbf{W}'_2\| \\
 & \leq \sqrt{2}(1 - \alpha_2)c_W P C_2 \left( \frac{\beta_2 B_2}{C_2} \right)^{\tilde{\alpha}} \|\text{vec}[\mathbf{W}_2] - \text{vec}[\mathbf{W}'_2]\|_2^{\tilde{\alpha}} \\
 & \quad + (1 - \alpha_2) \left( 2c_W^2 \beta_2 B_2 + \sqrt{2}(1 - \beta_2)c_W \right) \|\text{vec}[\mathbf{W}_2] - \text{vec}[\mathbf{W}'_2]\|_2.
 \end{aligned} \tag{52}$$

Besides,

$$\begin{aligned}
 & \|\delta_{ij}(\mathbf{W}_1) - \delta_{ij}(\mathbf{W}'_1)\|_2 \\
 & \leq (1 - \alpha_2) \left\| \sigma'(\mathbf{H}_{j^*}^{(1)}) \odot ((\delta_i(\mathbf{W}_1) - \delta_i(\mathbf{W}'_1))((1 - \beta_2)\mathbf{W}_2^\top + \beta_2\mathbf{I})) \right\|_2 \\
 & \quad + (1 - \alpha_2) \left\| (\sigma'(\mathbf{H}_{j^*}^{(1)})(\mathbf{W}_1) - \sigma'(\mathbf{H}_{j^*}^{(1)})(\mathbf{W}'_1)) \odot (\delta_i((1 - \beta_2)\mathbf{W}_2^\top + \beta_2\mathbf{I})) \right\|_2 \\
 & \leq (1 - \alpha_2) C_2 \|\delta_i(\mathbf{W}_1) - \delta_i(\mathbf{W}'_1)\|_2 + (1 - \alpha_2) C_2 \|\delta_i\| \left\| \sigma'(\mathbf{H}_{j^*}^{(1)})(\mathbf{W}_1) - \sigma'(\mathbf{H}_{j^*}^{(1)})(\mathbf{W}'_1) \right\|_2 \\
 & \leq (1 - \alpha_2) C_2 \|\delta_i(\mathbf{W}_1) - \delta_i(\mathbf{W}'_1)\|_2 + \sqrt{2}(1 - \alpha_2) c_W C_2 P \left\| \mathbf{H}_{j^*}^{(1)}(\mathbf{W}_1) - \mathbf{H}_{j^*}^{(1)}(\mathbf{W}'_1) \right\|_2^{\tilde{\alpha}} \\
 & \leq \sqrt{2} c_W C_2 P \left[ (1 - \alpha_2)^{1 + \tilde{\alpha}} C_2^{\tilde{\alpha}} \left\| g(\tilde{\mathbf{A}}) \right\|_\infty^{\tilde{\alpha}} + (1 - \alpha_2) \right] \left( \frac{\beta_1 B_1}{C_1} \right)^{\tilde{\alpha}} \|\text{vec}[\mathbf{W}_1] - \text{vec}[\mathbf{W}'_1]\|_2^{\tilde{\alpha}} \\
 & \quad + 2(1 - \alpha_2)^2 \beta_1 c_W^2 \left\| g(\tilde{\mathbf{A}}) \right\|_\infty \frac{B_1 C_2^2}{C_1} \|\text{vec}[\mathbf{W}_1] - \text{vec}[\mathbf{W}'_1]\|_2.
 \end{aligned} \tag{53}$$

Finally,

$$\begin{aligned}
 & \|\delta_{ij}(\mathbf{W}_0) - \delta_{ij}(\mathbf{W}'_0)\|_2 \\
 & \leq (1 - \alpha_2) \left\| \sigma'(\mathbf{H}_{j^*}^{(1)}) \odot ((\delta_i(\mathbf{W}_0) - \delta_i(\mathbf{W}'_0))((1 - \beta_2)\mathbf{W}_2^\top + \beta_2\mathbf{I})) \right\|_2 \\
 & \quad + (1 - \alpha_2) \left\| (\sigma'(\mathbf{H}_{j^*}^{(1)})(\mathbf{W}_0) - \sigma'(\mathbf{H}_{j^*}^{(1)})(\mathbf{W}'_0)) \odot (\delta_i((1 - \beta_2)\mathbf{W}_2^\top + \beta_2\mathbf{I})) \right\|_2 \\
 & \leq (1 - \alpha_2) C_2 \|\delta_i(\mathbf{W}_0) - \delta_i(\mathbf{W}'_0)\|_2 + \sqrt{2}(1 - \alpha_2) c_W C_2 P \left\| \mathbf{H}_{j^*}^{(1)}(\mathbf{W}_0) - \mathbf{H}_{j^*}^{(1)}(\mathbf{W}'_0) \right\|_2^{\tilde{\alpha}} \\
 & \leq \sqrt{2}(1 - \alpha_2) C_2 P c_W \left[ \left( \frac{B_2}{c_W} \right)^{\tilde{\alpha}} + \left( \frac{B_1}{c_W} \right)^{\tilde{\alpha}} \right] \|\text{vec}[\mathbf{W}_0] - \text{vec}[\mathbf{W}'_0]\|_2^{\tilde{\alpha}} \\
 & \quad + 2(1 - \alpha_2) c_W B_2 C_2 \|\text{vec}[\mathbf{W}_0] - \text{vec}[\mathbf{W}'_0]\|_2.
 \end{aligned} \tag{54}$$

Now we are ready to discuss each gradient term in the following four parts.

**Part F.** First

$$\begin{aligned}
 & \left\| \frac{\partial \ell(\mathbf{W}_0, \mathbf{W}_1, \mathbf{W}_2, \mathbf{W}_3; z_i)}{\partial \text{vec}[\mathbf{W}_3]}(\mathbf{W}_3) - \frac{\partial \ell(\mathbf{W}_0, \mathbf{W}_1, \mathbf{W}_2, \mathbf{W}_3; z_i)}{\partial \text{vec}[\mathbf{W}_3]}(\mathbf{W}'_3) \right\|_2 \\
 & \leq \max_{i \in [n]} \left\| \mathbf{H}_{i^*}^{(2)} \right\|_2 \|\hat{\mathbf{y}}_i(\mathbf{W}_3) - \hat{\mathbf{y}}_i(\mathbf{W}'_3)\|_2 \leq 2B_2^2 \|\text{vec}[\mathbf{W}_2] - \text{vec}[\mathbf{W}'_2]\|_2.
 \end{aligned}$$

Similarly,

$$\begin{aligned}
 & \left\| \frac{\partial \ell(\mathbf{W}_0, \mathbf{W}_1, \mathbf{W}_2, \mathbf{W}_3; z_i)}{\partial \text{vec}[\mathbf{W}_3]}(\mathbf{W}_2) - \frac{\partial \ell(\mathbf{W}_0, \mathbf{W}_1, \mathbf{W}_2, \mathbf{W}_3; z_i)}{\partial \text{vec}[\mathbf{W}_3]}(\mathbf{W}'_2) \right\|_2 \\
 & \leq \left\| \mathbf{H}_{i^*}^{(2)}(\mathbf{W}_2) - \mathbf{H}_{i^*}^{(2)}(\mathbf{W}'_2) \right\|_2 \|\hat{\mathbf{y}}_i - \mathbf{y}_i\|_2 + \max_{i \in [n]} \left\| \mathbf{H}_{i^*}^{(2)} \right\|_2 \|\hat{\mathbf{y}}_i(\mathbf{W}_2) - \hat{\mathbf{y}}_i(\mathbf{W}'_2)\|_2 \\
 & = (\sqrt{2} + 2c_W B_2) \Delta_{22} = (\sqrt{2} + 2c_W B_2) \frac{\beta_2 B_2}{C_2} \|\text{vec}[\mathbf{W}_2] - \text{vec}[\mathbf{W}'_2]\|_2.
 \end{aligned}$$

Similarly, we can obtain

$$\begin{aligned}
 & \left\| \frac{\partial \ell(\mathbf{W}_0, \mathbf{W}_1, \mathbf{W}_2, \mathbf{W}_3; z_i)}{\partial \text{vec}[\mathbf{W}_3]}(\mathbf{W}_1) - \frac{\partial \ell(\mathbf{W}_0, \mathbf{W}_1, \mathbf{W}_2, \mathbf{W}_3; z_i)}{\partial \text{vec}[\mathbf{W}_3]}(\mathbf{W}'_1) \right\|_2 \\
 & \leq \left\| \mathbf{H}_{i^*}^{(2)}(\mathbf{W}_1) - \mathbf{H}_{i^*}^{(2)}(\mathbf{W}'_1) \right\|_2 \|\hat{\mathbf{y}}_i - \mathbf{y}_i\|_2 + \max_{i \in [n]} \left\| \mathbf{H}_{i^*}^{(2)} \right\|_2 \|\hat{\mathbf{y}}_i(\mathbf{W}_1) - \hat{\mathbf{y}}_i(\mathbf{W}'_1)\|_2 \\
 & = (\sqrt{2} + 2c_W B_2) \Delta_{21} = (1 - \alpha_2) \beta_1 (\sqrt{2} + 2c_W B_2) \frac{B_1 C_2}{C_1} \left\| g(\tilde{\mathbf{A}}) \right\|_\infty \|\text{vec}[\mathbf{W}_1] - \text{vec}[\mathbf{W}'_1]\|_2,
 \end{aligned}$$

and

$$\begin{aligned}
 & \left\| \frac{\partial \ell(\mathbf{W}_0, \mathbf{W}_1, \mathbf{W}_2, \mathbf{W}_3; z_i)}{\partial \text{vec}[\mathbf{W}_3]}(\mathbf{W}_0) - \frac{\partial \ell(\mathbf{W}_0, \mathbf{W}_1, \mathbf{W}_2, \mathbf{W}_3; z_i)}{\partial \text{vec}[\mathbf{W}_3]}(\mathbf{W}'_0) \right\|_2 \\
 & \leq \left\| \mathbf{H}_{i*}^{(2)}(\mathbf{W}_0) - \mathbf{H}_{i*}^{(2)}(\mathbf{W}'_0) \right\|_2 \|\hat{\mathbf{y}}_i - \mathbf{y}_i\|_2 + \max_{i \in [n]} \left\| \mathbf{H}_{i*}^{(2)} \right\|_2 \|\hat{\mathbf{y}}_i(\mathbf{W}_0) - \hat{\mathbf{y}}_i(\mathbf{W}'_0)\|_2 \\
 & = (\sqrt{2} + 2c_W B_2) \Delta_{20} = (\sqrt{2} + 2c_W B_2) \frac{B_2}{c_W} \|\text{vec}[\mathbf{W}_0] - \text{vec}[\mathbf{W}'_0]\|_2.
 \end{aligned}$$

Denote by

$$\begin{aligned}
 P_{33} &= 2B_2^2, \quad P_{32} = (\sqrt{2} + 2c_W B_2) \frac{\beta_2 B_2}{C_2}, \quad P_{31} = (1 - \alpha_2) \beta_1 (\sqrt{2} + 2c_W B_2) \frac{B_1 C_2}{C_1} \left\| g(\tilde{\mathbf{A}}) \right\|_\infty, \\
 P_{30} &= (\sqrt{2} + 2c_W B_2) \frac{B_2}{c_W}, \quad \tilde{P}_{33} = \tilde{P}_{32} = \tilde{P}_{31} = \tilde{P}_{30} = 0,
 \end{aligned}$$

we obtain that  $\left\| \frac{\partial \ell(\mathbf{w}; z_i)}{\partial \text{vec}[\mathbf{W}_3]} - \frac{\partial \ell(\mathbf{w}'; z_i)}{\partial \text{vec}[\mathbf{W}_3]} \right\|_2 \leq \sum_{i=1}^4 P_{3i} \|\text{vec}[\mathbf{W}_i] - \text{vec}[\mathbf{W}'_i]\|_2 + \sum_{i=1}^4 \tilde{P}_{3i} \|\text{vec}[\mathbf{W}_i] - \text{vec}[\mathbf{W}'_i]\|_2^{\tilde{\alpha}}$ .

**Part G.** First,

$$\begin{aligned}
 & \left\| \frac{\partial \ell(\mathbf{W}_0, \mathbf{W}_1, \mathbf{W}_2, \mathbf{W}_3; z_i)}{\partial \text{vec}[\mathbf{W}_2]}(\mathbf{W}_3) - \frac{\partial \ell(\mathbf{W}_0, \mathbf{W}_1, \mathbf{W}_2, \mathbf{W}_3; z_i)}{\partial \text{vec}[\mathbf{W}_2]}(\mathbf{W}'_3) \right\|_2 \\
 & = \left\| \alpha_2 \beta_2 \mathbf{H}_{i*}^{(0)} + (1 - \alpha_2) \beta_2 \sum_{j=1}^n \left[ g(\tilde{\mathbf{A}}) \right]_{ij} \mathbf{H}_{j*}^{(1)} \right\|_2 \|\delta_i(\mathbf{W}_3) - \delta_i(\mathbf{W}'_3)\|_2 \\
 & \leq (\sqrt{2} + 2c_W B_2) \left( \alpha_2 \beta_2 \max_{i \in [n]} \left\| \mathbf{H}_{i*}^{(0)} \right\|_2 + (1 - \alpha_2) \beta_2 \left\| g(\tilde{\mathbf{A}}) \right\|_\infty \max_{i \in [n]} \left\| \mathbf{H}_{i*}^{(1)} \right\|_2 \right) \|\text{vec}[\mathbf{W}_3] - \text{vec}[\mathbf{W}'_3]\|_2 \\
 & \leq (\sqrt{2} + 2c_W B_2) \left( \alpha_2 \beta_2 c_X c_W + (1 - \alpha_2) \beta_2 B_1 \left\| g(\tilde{\mathbf{A}}) \right\|_\infty \right) \|\text{vec}[\mathbf{W}_3] - \text{vec}[\mathbf{W}'_3]\|_2 \\
 & = (\sqrt{2} + 2c_W B_2) \frac{\beta_2 B_2}{C_2} \|\text{vec}[\mathbf{W}_3] - \text{vec}[\mathbf{W}'_3]\|_2.
 \end{aligned}$$

Similarly,

$$\begin{aligned}
 & \left\| \frac{\partial \ell(\mathbf{W}_0, \mathbf{W}_1, \mathbf{W}_2, \mathbf{W}_3; z_i)}{\partial \text{vec}[\mathbf{W}_2]}(\mathbf{W}_2) - \frac{\partial \ell(\mathbf{W}_0, \mathbf{W}_1, \mathbf{W}_2, \mathbf{W}_3; z_i)}{\partial \text{vec}[\mathbf{W}_2]}(\mathbf{W}'_2) \right\|_2 \\
 & = \left\| \alpha_2 \beta_2 \mathbf{H}_{i*}^{(0)} + (1 - \alpha_2) \beta_2 \sum_{j=1}^n \left[ g(\tilde{\mathbf{A}}) \right]_{ij} \mathbf{H}_{j*}^{(1)} \right\|_2 \|\delta_i(\mathbf{W}_2) - \delta_i(\mathbf{W}'_2)\|_2 \\
 & \leq \sqrt{2} P_{c_W} \left( \frac{\beta_2 B_2}{C_2} \right)^{1+\tilde{\alpha}} \|\text{vec}[\mathbf{W}_2] - \text{vec}[\mathbf{W}'_2]\|_2^{\tilde{\alpha}} + c_W^2 \left( \frac{\beta_2 B_2}{C_2} \right)^2 \|\text{vec}[\mathbf{W}_2] - \text{vec}[\mathbf{W}'_2]\|_2.
 \end{aligned}$$

Besides,

$$\begin{aligned}
 & \left\| \frac{\partial \ell(\mathbf{W}_0, \mathbf{W}_1, \mathbf{W}_2, \mathbf{W}_3; z_i)}{\partial \text{vec}[\mathbf{W}_2]}(\mathbf{W}_1) - \frac{\partial \ell(\mathbf{W}_0, \mathbf{W}_1, \mathbf{W}_2, \mathbf{W}_3; z_i)}{\partial \text{vec}[\mathbf{W}_2]}(\mathbf{W}'_1) \right\|_2 \\
 & \leq \left\| \alpha_2 \beta_2 \mathbf{H}_{i*}^{(0)} + (1 - \alpha_2) \beta_2 \sum_{j=1}^n \left[ g(\tilde{\mathbf{A}}) \right]_{ij} \mathbf{H}_{j*}^{(1)} \right\|_2 \|\delta_i(\mathbf{W}_1) - \delta_i(\mathbf{W}'_1)\|_2 \\
 & \quad + \left\| (1 - \alpha_2) \beta_2 \sum_{j=1}^n \left[ g(\tilde{\mathbf{A}}) \right]_{ij} (\mathbf{H}_{j*}^{(1)}(\mathbf{W}_1) - \mathbf{H}_{j*}^{(1)}(\mathbf{W}'_1)) \right\|_2 \|\delta_i\|_2 \\
 & \leq \sqrt{2} P_{c_W} (1 - \alpha_2)^{\tilde{\alpha}} \beta_1^{\tilde{\alpha}} \left( \frac{B_1 C_2}{C_1} \right)^{\tilde{\alpha}} \left\| g(\tilde{\mathbf{A}}) \right\|_\infty^{\tilde{\alpha}} \frac{\beta_2 B_2}{C_2} \|\text{vec}[\mathbf{W}_1] - \text{vec}[\mathbf{W}'_1]\|_2^{\tilde{\alpha}} \\
 & \quad + (1 - \alpha_2) c_W \frac{B_1 \beta_1 \beta_2}{C_1} \left\| g(\tilde{\mathbf{A}}) \right\|_\infty \left[ c_W B_2 + \sqrt{2} \right] \|\text{vec}[\mathbf{W}_1] - \text{vec}[\mathbf{W}'_1]\|_2.
 \end{aligned}$$

Finally,

$$\begin{aligned}
 & \left\| \frac{\partial \ell(\mathbf{W}_0, \mathbf{W}_1, \mathbf{W}_2, \mathbf{W}_3; z_i)}{\partial \text{vec}[\mathbf{W}_2]}(\mathbf{W}_0) - \frac{\partial \ell(\mathbf{W}_0, \mathbf{W}_0, \mathbf{W}_2, \mathbf{W}_3; z_i)}{\partial \text{vec}[\mathbf{W}_2]}(\mathbf{W}'_0) \right\|_2 \\
 & \leq \left\| \alpha_2 \beta_2 \mathbf{H}_{i^*}^{(0)} + (1 - \alpha_2) \beta_2 \sum_{j=1}^n \left[ g(\tilde{\mathbf{A}}) \right]_{ij} \mathbf{H}_{j^*}^{(1)} \right\|_2 \|\delta_i(\mathbf{W}_0) - \delta_i(\mathbf{W}'_0)\|_2 \\
 & \quad + \left\| \alpha_2 \beta_2 (\mathbf{H}_{i^*}^{(0)}(\mathbf{W}_0) - \mathbf{H}_{i^*}^{(0)}(\mathbf{W}_0)) \right\|_2 \max_{i \in [n]} \|\delta_i\|_2 + (1 - \alpha_2) \beta_2 \left\| g(\tilde{\mathbf{A}}) \right\|_\infty \max_{i \in [n]} \|\delta_i\|_2 \left\| \mathbf{H}_{i^*}^{(1)}(\mathbf{W}_0) - \mathbf{H}_{i^*}^{(1)}(\mathbf{W}_0) \right\|_2 \\
 & \leq \sqrt{2} P_{cW} \left( \frac{B_2}{c_W} \right)^{\tilde{\alpha}} \frac{\beta_2 B_2}{C_2} \|\text{vec}[\mathbf{W}_0] - \text{vec}[\mathbf{W}'_0]\|_2^{\tilde{\alpha}} + \frac{\beta_2 B_2 (\sqrt{2} + B_2 c_W)}{C_2} \|\text{vec}[\mathbf{W}_0] - \text{vec}[\mathbf{W}'_0]\|_2.
 \end{aligned}$$

Denote by

$$\begin{aligned}
 P_{23} &= (\sqrt{2} + 2c_W B_2) \frac{\beta_2 B_2}{C_2}, \quad P_{22} = c_W^2 \left( \frac{\beta_2 B_2}{C_2} \right)^2, \\
 P_{21} &= (1 - \alpha_2) c_W \frac{B_1 \beta_1 \beta_2}{C_1} \left\| g(\tilde{\mathbf{A}}) \right\|_\infty [c_W B_2 + \sqrt{2}], \quad P_{20} = \frac{\beta_2 B_2 (\sqrt{2} + B_2 c_W)}{C_2} \\
 \tilde{P}_{23} &= 0, \quad \tilde{P}_{22} = \sqrt{2} P_{cW} \left( \frac{\beta_2 B_2}{C_2} \right)^{1+\tilde{\alpha}}, \quad \tilde{P}_{21} = \sqrt{2} P_{cW} (1 - \alpha_2)^{\tilde{\alpha}} \beta_1^{\tilde{\alpha}} \left( \frac{B_1 C_2}{C_1} \right)^{\tilde{\alpha}} \left\| g(\tilde{\mathbf{A}}) \right\|_\infty^{\tilde{\alpha}} \frac{\beta_2 B_2}{C_2}, \\
 \tilde{P}_{20} &= \sqrt{2} P_{cW} \left( \frac{B_2}{c_W} \right)^{\tilde{\alpha}} \frac{\beta_2 B_2}{C_2},
 \end{aligned}$$

we obtain that  $\left\| \frac{\partial \ell(\mathbf{w}; z_i)}{\partial \text{vec}[\mathbf{W}_2]} - \frac{\partial \ell(\mathbf{w}'; z_i)}{\partial \text{vec}[\mathbf{W}_2]} \right\|_2 \leq \sum_{i=1}^4 P_{2i} \|\text{vec}[\mathbf{W}_i] - \text{vec}[\mathbf{W}'_i]\|_2 + \sum_{i=1}^4 \tilde{P}_{2i} \|\text{vec}[\mathbf{W}_i] - \text{vec}[\mathbf{W}'_i]\|_2^{\tilde{\alpha}}$ .

**Part H.** Since (1)  $\mathbf{H}^{(0)}$  is variable related to  $\mathbf{W}_0$ ; (2)  $\delta_{ij}$  is variable related to  $\mathbf{W}_0, \mathbf{W}_1, \mathbf{W}_2, \mathbf{W}_3$ , we have

$$\begin{aligned}
 & \left\| \frac{\partial \ell(\mathbf{W}_0, \mathbf{W}_1, \mathbf{W}_2, \mathbf{W}_3; z_i)}{\partial \text{vec}[\mathbf{W}_1]}(\mathbf{W}_3) - \frac{\partial \ell(\mathbf{W}_0, \mathbf{W}_1, \mathbf{W}_2, \mathbf{W}_3; z_i)}{\partial \text{vec}[\mathbf{W}_1]}(\mathbf{W}'_3) \right\|_2 \\
 & = \alpha_1 \beta_1 \sum_{j=1}^n \left[ g(\tilde{\mathbf{A}}) \right]_{ij} \left\| \mathbf{H}_{j^*}^{(0)} \right\|_2 \|\delta_{ij}(\mathbf{W}_3) - \delta_{ij}(\mathbf{W}'_3)\|_2 \\
 & \quad + (1 - \alpha_1) \beta_1 \sum_{j=1}^n \sum_{k=1}^n \left[ g(\tilde{\mathbf{A}}) \right]_{ij} \left[ g(\tilde{\mathbf{A}}) \right]_{jk} \left\| \mathbf{H}_{k^*}^{(0)} \right\|_2 \|\delta_{ij}(\mathbf{W}_3) - \delta_{ij}(\mathbf{W}'_3)\|_2 \\
 & \leq c_X c_W (1 - \alpha_2) C_2 \left( \alpha_1 \beta_1 \left\| g(\tilde{\mathbf{A}}) \right\|_\infty + (1 - \alpha_1) \beta_1 \left\| g(\tilde{\mathbf{A}}) \right\|_\infty^2 \right) (\sqrt{2} + 2c_W B_2) \|\text{vec}[\mathbf{W}_3] - \text{vec}[\mathbf{W}'_3]\|_2 \\
 & = \frac{\beta_1 B_1 C_2}{C_1} \left\| g(\tilde{\mathbf{A}}) \right\|_\infty (\sqrt{2} + 2c_W B_2) \|\text{vec}[\mathbf{W}_3] - \text{vec}[\mathbf{W}'_3]\|_2,
 \end{aligned}$$

where we have used

$$\alpha_1 \beta_1 c_X c_W \left\| g(\tilde{\mathbf{A}}) \right\|_\infty + (1 - \alpha_1) \beta_1 c_X c_W \left\| g(\tilde{\mathbf{A}}) \right\|_\infty^2 = \beta_1 \frac{B_1}{C_1} \left\| g(\tilde{\mathbf{A}}) \right\|_\infty.$$

Similarly, we have

$$\begin{aligned}
 & \left\| \frac{\partial \ell(\mathbf{W}_0, \mathbf{W}_1, \mathbf{W}_2, \mathbf{W}_3; z_i)}{\partial \text{vec}[\mathbf{W}_1]}(\mathbf{W}_2) - \frac{\partial \ell(\mathbf{W}_0, \mathbf{W}_1, \mathbf{W}_2, \mathbf{W}_3; z_i)}{\partial \text{vec}[\mathbf{W}_1]}(\mathbf{W}'_2) \right\|_2 \\
 &= \alpha_1 \beta_1 \sum_{j=1}^n \left[ g(\tilde{\mathbf{A}}) \right]_{ij} \left\| \mathbf{H}_{j*}^{(0)} \right\|_2 \left\| \delta_{ij}(\mathbf{W}_2) - \delta_{ij}(\mathbf{W}'_2) \right\|_2 \\
 & \quad + (1 - \alpha_1) \beta_1 \sum_{j=1}^n \sum_{k=1}^n \left[ g(\tilde{\mathbf{A}}) \right]_{ij} \left[ g(\tilde{\mathbf{A}}) \right]_{jk} \left\| \mathbf{H}_{k*}^{(0)} \right\|_2 \left\| \delta_{ij}(\mathbf{W}_2) - \delta_{ij}(\mathbf{W}'_2) \right\|_2 \\
 & \leq \sqrt{2} (1 - \alpha_2) c_W P \frac{\beta_1 B_1 C_2}{C_1} \left\| g(\tilde{\mathbf{A}}) \right\|_\infty \left( \frac{\beta_2 B_2}{C_2} \right)^{\tilde{\alpha}} \left\| \text{vec}[\mathbf{W}_2] - \text{vec}[\mathbf{W}'_2] \right\|_2^{\tilde{\alpha}} \\
 & \quad + (1 - \alpha_2) \frac{\beta_1 B_1}{C_1} \left\| g(\tilde{\mathbf{A}}) \right\|_\infty \left( c_W^2 \beta_2 B_2 + \sqrt{2} (1 - \beta_2) c_W \right) \left\| \text{vec}[\mathbf{W}_2] - \text{vec}[\mathbf{W}'_2] \right\|_2.
 \end{aligned}$$

Besides,

$$\begin{aligned}
 & \left\| \frac{\partial \ell(\mathbf{W}_0, \mathbf{W}_1, \mathbf{W}_2, \mathbf{W}_3; z_i)}{\partial \text{vec}[\mathbf{W}_1]}(\mathbf{W}_1) - \frac{\partial \ell(\mathbf{W}_0, \mathbf{W}_1, \mathbf{W}_2, \mathbf{W}_3; z_i)}{\partial \text{vec}[\mathbf{W}_1]}(\mathbf{W}'_1) \right\|_2 \\
 &= \alpha_1 \beta_1 \sum_{j=1}^n \left[ g(\tilde{\mathbf{A}}) \right]_{ij} \left\| \mathbf{H}_{j*}^{(0)} \right\|_2 \left\| \delta_{ij}(\mathbf{W}_1) - \delta_{ij}(\mathbf{W}'_1) \right\|_2 \\
 & \quad + (1 - \alpha_1) \beta_1 \sum_{j=1}^n \sum_{k=1}^n \left[ g(\tilde{\mathbf{A}}) \right]_{ij} \left[ g(\tilde{\mathbf{A}}) \right]_{jk} \left\| \mathbf{H}_{k*}^{(0)} \right\|_2 \left\| \delta_{ij}(\mathbf{W}_1) - \delta_{ij}(\mathbf{W}'_1) \right\|_2 \\
 & \leq \sqrt{2} c_W C_2 P \left[ (1 - \alpha_2)^{1+\tilde{\alpha}} C_2^{\tilde{\alpha}} \left\| g(\tilde{\mathbf{A}}) \right\|_\infty^{\tilde{\alpha}} + (1 - \alpha_2) \right] \left( \frac{\beta_1 B_1}{C_1} \right)^{\tilde{\alpha}+1} \left\| g(\tilde{\mathbf{A}}) \right\|_\infty \left\| \text{vec}[\mathbf{W}_1] - \text{vec}[\mathbf{W}'_1] \right\|_2^{\tilde{\alpha}} \\
 & \quad + (1 - \alpha_2)^2 \beta_1^2 c_W^2 \frac{B_1^2 C_2^2}{C_1^2} \left\| g(\tilde{\mathbf{A}}) \right\|_\infty \left\| \text{vec}[\mathbf{W}_1] - \text{vec}[\mathbf{W}'_1] \right\|_2.
 \end{aligned}$$

Finally,

$$\begin{aligned}
 & \left\| \frac{\partial \ell(\mathbf{W}_0, \mathbf{W}_1, \mathbf{W}_2, \mathbf{W}_3; z_i)}{\partial \text{vec}[\mathbf{W}_1]}(\mathbf{W}_0) - \frac{\partial \ell(\mathbf{W}_0, \mathbf{W}_1, \mathbf{W}_2, \mathbf{W}_3; z_i)}{\partial \text{vec}[\mathbf{W}_1]}(\mathbf{W}'_0) \right\|_2 \\
 &= \alpha_1 \beta_1 \sum_{j=1}^n \left[ g(\tilde{\mathbf{A}}) \right]_{ij} \left\| \mathbf{H}_{j*}^{(0)} \right\|_2 \left\| \delta_{ij}(\mathbf{W}_0) - \delta_{ij}(\mathbf{W}'_0) \right\|_2 \\
 & \quad + (1 - \alpha_1) \beta_1 \sum_{j=1}^n \sum_{k=1}^n \left[ g(\tilde{\mathbf{A}}) \right]_{ij} \left[ g(\tilde{\mathbf{A}}) \right]_{jk} \left\| \mathbf{H}_{k*}^{(0)} \right\|_2 \left\| \delta_{ij}(\mathbf{W}_0) - \delta_{ij}(\mathbf{W}'_0) \right\|_2 \\
 & \leq \sqrt{2} (1 - \alpha_2) c_W P \left[ \left( \frac{B_2}{c_W} \right)^{\tilde{\alpha}} + \left( \frac{B_1}{c_W} \right)^{\tilde{\alpha}} \right] \frac{\beta_1 B_1 C_2}{C_1} \left\| g(\tilde{\mathbf{A}}) \right\|_\infty \left\| \text{vec}[\mathbf{W}_0] - \text{vec}[\mathbf{W}'_0] \right\|_2^{\tilde{\alpha}} \\
 & \quad + (1 - \alpha_2) c_W \frac{\beta_1 B_1 B_2}{C_1} \left\| g(\tilde{\mathbf{A}}) \right\|_\infty \left\| \text{vec}[\mathbf{W}_0] - \text{vec}[\mathbf{W}'_0] \right\|_2.
 \end{aligned}$$

Denote by

$$\begin{aligned}
 P_{13} &= \frac{\beta_1 B_1 C_2}{C_1} \left\| g(\tilde{\mathbf{A}}) \right\|_{\infty} (\sqrt{2} + 2c_W B_2), \quad P_{12} = (1 - \alpha_2) \frac{\beta_1 B_1}{C_1} \left\| g(\tilde{\mathbf{A}}) \right\|_{\infty} \left( c_W^2 \beta_2 B_2 + \sqrt{2}(1 - \beta_2)c_W \right) \\
 P_{11} &= (1 - \alpha_2)^2 \beta_1^2 c_W^2 \frac{B_1^2 C_2^2}{C_1^2} \left\| g(\tilde{\mathbf{A}}) \right\|_{\infty}, \quad P_{10} = (1 - \alpha_2)c_W \frac{\beta_1 B_1 B_2}{C_1} \left\| g(\tilde{\mathbf{A}}) \right\|_{\infty} \\
 \tilde{P}_{13} &= 0, \quad \tilde{P}_{12} = \sqrt{2}(1 - \alpha_2)c_W P \frac{\beta_1 B_1 C_2}{C_1} \left\| g(\tilde{\mathbf{A}}) \right\|_{\infty} \left( \frac{\beta_2 B_2}{C_2} \right)^{\tilde{\alpha}} \\
 \tilde{P}_{11} &= \sqrt{2}c_W C_2 P \left[ (1 - \alpha_2)^{1+\tilde{\alpha}} C_2^{\tilde{\alpha}} \left\| g(\tilde{\mathbf{A}}) \right\|_{\infty}^{\tilde{\alpha}} + (1 - \alpha_2) \right] \left( \frac{\beta_1 B_1}{C_1} \right)^{\tilde{\alpha}+1} \left\| g(\tilde{\mathbf{A}}) \right\|_{\infty}, \\
 P_{10} &= \sqrt{2}(1 - \alpha_2)c_W P \left[ \left( \frac{B_2}{c_W} \right)^{\tilde{\alpha}} + \left( \frac{B_1}{c_W} \right)^{\tilde{\alpha}} \right] \frac{\beta_1 B_1 C_2}{C_1} \left\| g(\tilde{\mathbf{A}}) \right\|_{\infty},
 \end{aligned}$$

we obtain that  $\left\| \frac{\partial \ell(\mathbf{w}; z_i)}{\partial \text{vec}[\mathbf{W}_1]} - \frac{\partial \ell(\mathbf{w}'; z_i)}{\partial \text{vec}[\mathbf{W}_1]} \right\|_2 \leq \sum_{i=1}^4 P_{1i} \|\text{vec}[\mathbf{W}_i] - \text{vec}[\mathbf{W}'_i]\|_2 + \sum_{i=1}^4 \tilde{P}_{1i} \|\text{vec}[\mathbf{W}_i] - \text{vec}[\mathbf{W}'_i]\|_2^{\tilde{\alpha}}$ .

**Part I.** First we have

$$\begin{aligned}
 & \left\| \frac{\partial \ell(\mathbf{W}_0, \mathbf{W}_1, \mathbf{W}_2, \mathbf{W}_3; z_i)}{\partial \text{vec}[\mathbf{W}_0]}(\mathbf{W}_3) - \frac{\partial \ell(\mathbf{W}_0, \mathbf{W}_1, \mathbf{W}_2, \mathbf{W}_3; z_i)}{\partial \text{vec}[\mathbf{W}_0]}(\mathbf{W}'_3) \right\|_2 \\
 & \leq \alpha_2 \|\mathbf{X}_{i*}\|_2 \|\delta_i(\mathbf{W}_3) - \delta_i(\mathbf{W}'_3)\|_2 \|(1 - \beta_2)\mathbf{I} + \beta_2 \mathbf{W}_2^{\top}\|_2 \|\mathbf{H}_{i*}^{(0)}\|_2 \\
 & \quad + \alpha_1 \sum_{j=1}^n \left[ g(\tilde{\mathbf{A}}) \right]_{ij} \|\mathbf{X}_{j*}\|_2 \|\delta_{ij}(\mathbf{W}_3) - \delta_{ij}(\mathbf{W}'_3)\|_2 \|(1 - \beta_1)\mathbf{I} + \beta_1 \mathbf{W}_1^{\top}\|_2 \|\mathbf{H}_{j*}^{(0)}\|_2 \\
 & \quad + (1 - \alpha_1) \sum_{j=1}^n \sum_{k=1}^n \left[ g(\tilde{\mathbf{A}}) \right]_{ij} \left[ g(\tilde{\mathbf{A}}) \right]_{jk} \|\mathbf{X}_{k*}\|_2 \|\delta_{ij}(\mathbf{W}_3) - \delta_{ij}(\mathbf{W}'_3)\|_2 \|(1 - \beta_1)\mathbf{I} + \beta_1 \mathbf{W}_1^{\top}\|_2 \|\mathbf{H}_{k*}^{(0)}\|_2 \\
 & \leq \left( \alpha_2 C_2 + \alpha_1 (1 - \alpha_2) C_1 C_2 \left\| g(\tilde{\mathbf{A}}) \right\|_{\infty} + (1 - \alpha_1)(1 - \alpha_2) C_1 C_2 \left\| g(\tilde{\mathbf{A}}) \right\|_{\infty}^2 \right) c_X^2 c_W \\
 & \quad \times (\sqrt{2} + 2c_W B_2) \|\text{vec}[\mathbf{W}_3] - \text{vec}[\mathbf{W}'_3]\|_2 \\
 & = c_X B_2 (\sqrt{2} + 2c_W B_2) \|\text{vec}[\mathbf{W}_3] - \text{vec}[\mathbf{W}'_3]\|_2.
 \end{aligned}$$

Similarly,

$$\begin{aligned}
 & \left\| \frac{\partial \ell(\mathbf{W}_0, \mathbf{W}_1, \mathbf{W}_2, \mathbf{W}_3; z_i)}{\partial \text{vec}[\mathbf{W}_0]}(\mathbf{W}_2) - \frac{\partial \ell(\mathbf{W}_0, \mathbf{W}_1, \mathbf{W}_2, \mathbf{W}_3; z_i)}{\partial \text{vec}[\mathbf{W}_0]}(\mathbf{W}'_2) \right\|_2 \\
 & \leq \alpha_2 \|\mathbf{X}_{i*}\|_2 \|\delta_i(\mathbf{W}_2) - \delta_i(\mathbf{W}'_2)\|_2 \|(1 - \beta_2)\mathbf{I} + \beta_2 \mathbf{W}_2^{\top}\|_2 \|\mathbf{H}_{i*}^{(0)}\|_2 \\
 & \quad + \alpha_2 \beta_2 \|\mathbf{X}_{i*}\|_2 \|\delta_i\|_2 \|\mathbf{H}_{i*}^{(0)}\|_2 \|\mathbf{W}_2 - \mathbf{W}'_2\| \\
 & \quad + \alpha_1 \sum_{j=1}^n \left[ g(\tilde{\mathbf{A}}) \right]_{ij} \|\mathbf{X}_{j*}\|_2 \|\delta_{ij}(\mathbf{W}_2) - \delta_{ij}(\mathbf{W}'_2)\|_2 \|(1 - \beta_1)\mathbf{I} + \beta_1 \mathbf{W}_1^{\top}\|_2 \|\mathbf{H}_{j*}^{(0)}\|_2 \\
 & \quad + (1 - \alpha_1) \sum_{j=1}^n \sum_{k=1}^n \left[ g(\tilde{\mathbf{A}}) \right]_{ij} \left[ g(\tilde{\mathbf{A}}) \right]_{jk} \|\mathbf{X}_{k*}\|_2 \|\delta_{ij}(\mathbf{W}_2) - \delta_{ij}(\mathbf{W}'_2)\|_2 \|(1 - \beta_1)\mathbf{I} + \beta_1 \mathbf{W}_1^{\top}\|_2 \|\mathbf{H}_{k*}^{(0)}\|_2 \\
 & \leq B_2 c_X c_W \sqrt{2} P \left( \frac{\beta_2 B_2}{C_2} \right)^{\tilde{\alpha}} \|\text{vec}[\mathbf{W}_2] - \text{vec}[\mathbf{W}'_2]\|_2^{\tilde{\alpha}} \\
 & \quad + \left[ 2c_X c_W^2 \beta_2 \frac{B_2^2}{C_2} + \sqrt{2} \alpha_2 \beta_2 c_X^2 c_W^2 + \sqrt{2}(1 - \alpha_2)(1 - \beta_2) c_X c_W B_1 \left\| g(\tilde{\mathbf{A}}) \right\|_{\infty} \right] \|\text{vec}[\mathbf{W}_2] - \text{vec}[\mathbf{W}'_2]\|_2.
 \end{aligned}$$

Besides,

$$\begin{aligned}
 & \left\| \frac{\partial \ell(\mathbf{W}_0, \mathbf{W}_1, \mathbf{W}_2, \mathbf{W}_3; z_i)}{\partial \text{vec}[\mathbf{W}_0]}(\mathbf{W}_1) - \frac{\partial \ell(\mathbf{W}_0, \mathbf{W}_1, \mathbf{W}_2, \mathbf{W}_3; z_i)}{\partial \text{vec}[\mathbf{W}_0]}(\mathbf{W}'_1) \right\|_2 \\
 \leq & \alpha_2 \|\mathbf{X}_{i*}\|_2 \|\delta_i(\mathbf{W}_1) - \delta_i(\mathbf{W}'_1)\|_2 \|(1 - \beta_2)\mathbf{I} + \beta_2 \mathbf{W}_2^\top\|_2 \|\mathbf{H}_{i*}^{(0)}\|_2 \\
 & + \alpha_1 \beta_1 \sum_{j=1}^n \left[ g(\tilde{\mathbf{A}}) \right]_{ij} \|\mathbf{X}_{j*}\|_2 \|\delta_{ij}\|_2 \|\mathbf{H}_{j*}^{(0)}\|_2 \|\mathbf{W}_1 - \mathbf{W}'_1\| \\
 & + \alpha_1 \sum_{j=1}^n \left[ g(\tilde{\mathbf{A}}) \right]_{ij} \|\mathbf{X}_{j*}\|_2 \|\delta_{ij}(\mathbf{W}_1) - \delta_{ij}(\mathbf{W}'_1)\|_2 \|(1 - \beta_1)\mathbf{I} + \beta_1 \mathbf{W}_1^\top\|_2 \|\mathbf{H}_{j*}^{(0)}\|_2 \\
 & + (1 - \alpha_1) \sum_{j=1}^n \sum_{k=1}^n \left[ g(\tilde{\mathbf{A}}) \right]_{ij} \left[ g(\tilde{\mathbf{A}}) \right]_{jk} \|\mathbf{X}_{k*}\|_2 \|\delta_{ij}(\mathbf{W}_1) - \delta_{ij}(\mathbf{W}'_1)\|_2 \|(1 - \beta_1)\mathbf{I} + \beta_1 \mathbf{W}_1^\top\|_2 \|\mathbf{H}_{k*}^{(0)}\|_2 \\
 \leq & \left( B_2(1 - \alpha_2)^{\tilde{\alpha}} C_2^{\tilde{\alpha}} \left\| g(\tilde{\mathbf{A}}) \right\|_{\infty}^{\tilde{\alpha}} + B_1 C_2 \left\| g(\tilde{\mathbf{A}}) \right\|_{\infty} \right) \sqrt{2} P c_X c_W \left( \frac{\beta_1 B_1}{C_1} \right)^{\tilde{\alpha}} \|\text{vec}[\mathbf{W}_1] - \text{vec}[\mathbf{W}'_1]\|_2^{\tilde{\alpha}} \\
 & + \left( 2(1 - \alpha_2) \beta_1 c_X c_W^2 \frac{B_1 B_2 C_2}{C_1} + \sqrt{2} \alpha_1 (1 - \alpha_2) \beta_1 c_X^2 c_W^2 C_2 \right) \left\| g(\tilde{\mathbf{A}}) \right\|_{\infty} \|\text{vec}[\mathbf{W}_1] - \text{vec}[\mathbf{W}'_1]\|_2.
 \end{aligned}$$

Finally,

$$\begin{aligned}
 & \left\| \frac{\partial \ell(\mathbf{W}_0, \mathbf{W}_1, \mathbf{W}_2, \mathbf{W}_3; z_i)}{\partial \text{vec}[\mathbf{W}_0]}(\mathbf{W}_0) - \frac{\partial \ell(\mathbf{W}_0, \mathbf{W}_1, \mathbf{W}_2, \mathbf{W}_3; z_i)}{\partial \text{vec}[\mathbf{W}_0]}(\mathbf{W}'_0) \right\|_2 \\
 \leq & \alpha_2 \|\mathbf{X}_{i*}\|_2 \|\delta_i(\mathbf{W}_0) - \delta_i(\mathbf{W}'_0)\|_2 \|(1 - \beta_2)\mathbf{I} + \beta_2 \mathbf{W}_2^\top\|_2 \|\mathbf{H}_{i*}^{(0)}\|_2 \\
 & + \alpha_2 \|\mathbf{X}_{i*}\|_2 \|\delta_i((1 - \beta_2)\mathbf{I} + \beta_2 \mathbf{W}_2^\top)\|_2 \|\mathbf{H}_{i*}^{(0)}(\mathbf{W}_0) - \mathbf{H}_{i*}^{(0)}(\mathbf{W}'_0)\|_2 \\
 & + \alpha_1 \sum_{j=1}^n \left[ g(\tilde{\mathbf{A}}) \right]_{ij} \|\mathbf{X}_{j*}\|_2 \|\delta_{ij}(\mathbf{W}_1) - \delta_{ij}(\mathbf{W}'_1)\|_2 \|(1 - \beta_1)\mathbf{I} + \beta_1 \mathbf{W}_1^\top\|_2 \|\mathbf{H}_{j*}^{(0)}\|_2 \\
 & + \alpha_1 \sum_{j=1}^n \left[ g(\tilde{\mathbf{A}}) \right]_{ij} \|\mathbf{X}_{j*}\|_2 \|\delta_{ij}((1 - \beta_1)\mathbf{I} + \beta_1 \mathbf{W}_1^\top)\|_2 \|\mathbf{H}_{j*}^{(0)}(\mathbf{W}_0) - \mathbf{H}_{j*}^{(0)}(\mathbf{W}'_0)\|_2 \\
 & + (1 - \alpha_1) \sum_{j=1}^n \sum_{k=1}^n \left[ g(\tilde{\mathbf{A}}) \right]_{ij} \left[ g(\tilde{\mathbf{A}}) \right]_{jk} \|\mathbf{X}_{k*}\|_2 \|\delta_{ij}(\mathbf{W}_0) - \delta_{ij}(\mathbf{W}'_0)\|_2 \|(1 - \beta_1)\mathbf{I} + \beta_1 \mathbf{W}_1^\top\|_2 \|\mathbf{H}_{k*}^{(0)}\|_2 \\
 & + (1 - \alpha_1) \sum_{j=1}^n \sum_{k=1}^n \left[ g(\tilde{\mathbf{A}}) \right]_{ij} \left[ g(\tilde{\mathbf{A}}) \right]_{jk} \|\mathbf{X}_{k*}\|_2 \|\delta_{ij}((1 - \beta_1)\mathbf{I} + \beta_1 \mathbf{W}_1^\top)\|_2 \|\mathbf{H}_{k*}^{(0)}(\mathbf{W}_0) - \mathbf{H}_{k*}^{(0)}(\mathbf{W}'_0)\|_2 \\
 \leq & c_X B_2 (\sqrt{2} + 2c_W B_2) \|\text{vec}[\mathbf{W}_0] - \text{vec}[\mathbf{W}'_0]\| \\
 & + \sqrt{2} c_X c_W P \left[ B_2 \left( \frac{B_2}{c_W} \right)^{\tilde{\alpha}} + (1 - \alpha_2) C_2 \left\| g(\tilde{\mathbf{A}}) \right\|_{\infty} B_1 \left( \frac{B_1}{c_W} \right)^{\tilde{\alpha}} \right] \|\text{vec}[\mathbf{W}_0] - \text{vec}[\mathbf{W}'_0]\|_2^{\tilde{\alpha}}.
 \end{aligned}$$



Denote by

$$\begin{aligned}
 P_{03} &= c_X B_2 (\sqrt{2} + 2c_W B_2), \quad P_{02} = \sqrt{2} \alpha_2 \beta_2 c_X^2 c_W^2 + \sqrt{2} (1 - \alpha_2) (1 - \beta_2) c_X c_W B_1 \left\| g(\tilde{\mathbf{A}}) \right\|_\infty, \\
 P_{01} &= \left( 2(1 - \alpha_2) \beta_1 c_X c_W^2 \frac{B_1 B_2 C_2}{C_1} + \sqrt{2} \alpha_1 (1 - \alpha_2) \beta_1 c_X^2 c_W^2 C_2 \right) \left\| g(\tilde{\mathbf{A}}) \right\|_\infty \quad P_{00} = c_X B_2 (\sqrt{2} + 2c_W B_2), \\
 \tilde{P}_{03} &= 0, \quad \tilde{P}_{02} = B_2 c_X c_W \sqrt{2} P \left( \frac{\beta_2 B_2}{C_2} \right)^{\tilde{\alpha}}, \\
 \tilde{P}_{01} &= \left( B_2 (1 - \alpha_2)^{\tilde{\alpha}} C_2^{\tilde{\alpha}} \left\| g(\tilde{\mathbf{A}}) \right\|_\infty^{\tilde{\alpha}} + B_1 C_2 \left\| g(\tilde{\mathbf{A}}) \right\|_\infty \right) \sqrt{2} P c_X c_W \left( \frac{\beta_1 B_1}{C_1} \right)^{\tilde{\alpha}}, \\
 \tilde{P}_{00} &= \sqrt{2} c_X c_W P \left[ B_2 \left( \frac{B_2}{c_W} \right)^{\tilde{\alpha}} + (1 - \alpha_2) C_2 \left\| g(\tilde{\mathbf{A}}) \right\|_\infty B_1 \left( \frac{B_1}{c_W} \right)^{\tilde{\alpha}} \right],
 \end{aligned}$$

we obtain that  $\left\| \frac{\partial \ell(\mathbf{w}; z_i)}{\partial \text{vec}[\mathbf{W}_0]} - \frac{\partial \ell(\mathbf{w}'; z_i)}{\partial \text{vec}[\mathbf{W}_0]} \right\|_2 \leq \sum_{i=1}^4 P_{0i} \|\text{vec}[\mathbf{W}_i] - \text{vec}[\mathbf{W}'_i]\|_2 + \sum_{i=1}^4 \tilde{P}_{0i} \|\text{vec}[\mathbf{W}_i] - \text{vec}[\mathbf{W}'_i]\|_2^{\tilde{\alpha}}$ . Combining the results in Part F, Part G, Part H, Part I, we conclude that  $\|\nabla \ell(\mathbf{w}) - \nabla \ell(\mathbf{w}')\|_2 \leq P_{\mathcal{F}} \max\{\|\mathbf{w} - \mathbf{w}'\|_2, \|\mathbf{w} - \mathbf{w}'\|_2^{\tilde{\alpha}}\}$  holds where  $\mathbf{w} = [\text{vec}[\mathbf{W}_0]; \text{vec}[\mathbf{W}_1]; \text{vec}[\mathbf{W}_2]; \text{vec}[\mathbf{W}_3]]$  by Lemma A.4.

### B.2.3. PROOF OF PROPOSITION 4.13

For two layer SGC, we have  $g(\tilde{\mathbf{A}}) = \tilde{\mathbf{A}}^2$ . Note that

$$\begin{aligned}
 |\ell(\mathbf{W}_1 \mathbf{W}_2; z_i) - \ell(\mathbf{W}_1 \mathbf{W}'_2; z_i)| &\leq \sqrt{2} \left\| \sum_{j=1}^n \left[ g(\tilde{\mathbf{A}}) \right]_{ij} \mathbf{X}_{j*} \mathbf{W}_1 (\mathbf{W}_2 - \mathbf{W}'_2) \right\|_2 \\
 &\leq \sqrt{2} \sum_{j=1}^n \left[ g(\tilde{\mathbf{A}}) \right]_{ij} \|\mathbf{X}_{j*} \mathbf{W}_1\|_2 \|\mathbf{W}_2 - \mathbf{W}'_2\| \\
 &\leq c_X c_W \sqrt{2} \left\| g(\tilde{\mathbf{A}}) \right\|_\infty \|\mathbf{W}_2 - \mathbf{W}'_2\|.
 \end{aligned}$$

Similarly,

$$|\ell(\mathbf{W}_1 \mathbf{W}_2; z_i) - \ell(\mathbf{W}'_1 \mathbf{W}_2; z_i)| \leq c_X c_W \sqrt{2} \left\| g(\tilde{\mathbf{A}}) \right\|_\infty \|\mathbf{W}_1 - \mathbf{W}'_1\|.$$

Denote by  $L_1 = L_2 = c_X c_W \sqrt{2} \left\| g(\tilde{\mathbf{A}}) \right\|_\infty$ , we conclude that  $|\ell(\mathbf{w}) - \ell(\mathbf{w}')| \leq L_{\mathcal{F}} \|\mathbf{w} - \mathbf{w}'\|_2$  holds with  $L_{\mathcal{F}} = 2c_X c_W \left\| g(\tilde{\mathbf{A}}) \right\|_\infty$  by Lemma A.4. Then we discuss the Hölder smoothness. By the chain rule, the gradients are

$$\begin{aligned}
 \frac{\partial \ell(\mathbf{W}_1, \mathbf{W}_2; z_i)}{\partial \text{vec}[\mathbf{W}_2]} &= \sum_{j=1}^n \left[ g(\tilde{\mathbf{A}}) \right]_{ij} (\hat{\mathbf{y}}_i - \mathbf{y}_i) \otimes (\mathbf{X}_{j*} \mathbf{W}_1), \\
 \frac{\partial \ell(\mathbf{W}_1, \mathbf{W}_2; z_i)}{\partial \text{vec}[\mathbf{W}_1]} &= \sum_{j=1}^n \left[ g(\tilde{\mathbf{A}}) \right]_{ij} ((\hat{\mathbf{y}}_i - \mathbf{y}_i) \mathbf{W}_2^\top) \otimes \mathbf{X}_{j*}.
 \end{aligned}$$

First,

$$\begin{aligned}
 &\left\| \frac{\partial \ell(\mathbf{W}_1, \mathbf{W}_2; z_i)}{\partial \text{vec}[\mathbf{W}_2]}(\mathbf{W}_1) - \frac{\partial \ell(\mathbf{W}_1, \mathbf{W}_2; z_i)}{\partial \text{vec}[\mathbf{W}_2]}(\mathbf{W}'_1) \right\|_2 \\
 &\leq \sum_{j=1}^n \left[ g(\tilde{\mathbf{A}}) \right]_{ij} \|(\hat{\mathbf{y}}_i - \mathbf{y}_i) \otimes (\mathbf{X}_{j*} (\mathbf{W}_1 - \mathbf{W}'_1))\| + \sum_{j=1}^n \left[ g(\tilde{\mathbf{A}}) \right]_{ij} \|(\hat{\mathbf{y}}_i(\mathbf{W}_2) - \mathbf{y}_i(\mathbf{W}'_2)) \otimes (\mathbf{X}_{j*} \mathbf{W}_1)\| \\
 &\leq \left( \sqrt{2} c_X \left\| g(\tilde{\mathbf{A}}) \right\|_\infty + c_X^2 c_W^2 \left\| g(\tilde{\mathbf{A}}) \right\|_\infty^2 \right) \|\text{vec}[\mathbf{W}_1] - \text{vec}[\mathbf{W}'_1]\|.
 \end{aligned}$$

Similarly,

$$\begin{aligned} & \left\| \frac{\partial \ell(\mathbf{W}_1, \mathbf{W}_2; z_i)}{\partial \text{vec}[\mathbf{W}_2]}(\mathbf{W}_2) - \frac{\partial \ell(\mathbf{W}_1, \mathbf{W}_2; z_i)}{\partial \text{vec}[\mathbf{W}_2]}(\mathbf{W}'_2) \right\|_2 \\ & \leq \sum_{j=1}^n \left[ g(\tilde{\mathbf{A}}) \right]_{ij} \|\hat{\mathbf{y}}_i(\mathbf{W}_2) - \mathbf{y}_i(\mathbf{W}'_2)\| \otimes (\mathbf{X}_{j*} \mathbf{W}_1) \leq c_X^2 c_W^2 \left\| g(\tilde{\mathbf{A}}) \right\|_\infty^2 \|\text{vec}[\mathbf{W}_2] - \text{vec}[\mathbf{W}'_2]\|. \end{aligned}$$

Denote by  $P_{21} = \sqrt{2}c_X \left\| g(\tilde{\mathbf{A}}) \right\|_\infty + c_X^2 c_W^2 \left\| g(\tilde{\mathbf{A}}) \right\|_\infty^2$ ,  $P_{22} = c_X^2 c_W^2 \left\| g(\tilde{\mathbf{A}}) \right\|_\infty^2$  and  $\tilde{P}_{21} = \tilde{P}_{22} = 0$ , we obtain that  $\left\| \frac{\partial \ell(\mathbf{w}; z_i)}{\partial \text{vec}[\mathbf{W}_2]} - \frac{\partial \ell(\mathbf{w}'; z_i)}{\partial \text{vec}[\mathbf{W}_2]} \right\|_2 \leq \sum_{i=1}^2 P_{2i} \|\text{vec}[\mathbf{W}_i] - \text{vec}[\mathbf{W}'_i]\|_2 + \tilde{P}_{2i} \|\text{vec}[\mathbf{W}_i] - \text{vec}[\mathbf{W}'_i]\|_2^{\tilde{\alpha}}$ . By the same way, denote by  $P_{11} = P_{22}$ ,  $\tilde{P}_{11} = \tilde{P}_{22}$  and  $P_{12} = P_{21}$ ,  $\tilde{P}_{12} = \tilde{P}_{21}$  as well as  $\alpha_{11} = \alpha_{12} = 1$ , we obtain that  $\left\| \frac{\partial \ell(\mathbf{w}; z_i)}{\partial \text{vec}[\mathbf{W}_1]} - \frac{\partial \ell(\mathbf{w}'; z_i)}{\partial \text{vec}[\mathbf{W}_1]} \right\|_2 \leq \sum_{i=1}^2 P_{1i} \|\text{vec}[\mathbf{W}_i] - \text{vec}[\mathbf{W}'_i]\|_2 + \tilde{P}_{1i} \|\text{vec}[\mathbf{W}_i] - \text{vec}[\mathbf{W}'_i]\|_2^{\tilde{\alpha}}$ . By Lemma A.4, we conclude that  $\|\nabla \ell(\mathbf{w}) - \nabla \ell(\mathbf{w}')\|_2 \leq P_{\mathcal{F}} \max\{\|\mathbf{w} - \mathbf{w}'\|_2, \|\mathbf{w} - \mathbf{w}'\|_2^{\tilde{\alpha}}\}$  holds where  $\mathbf{w} = [\text{vec}[\mathbf{W}_1]; \text{vec}[\mathbf{W}_2]]$ .

#### B.2.4. PROOF OF PROPOSITION 4.14

We first show that the objective  $\ell(\mathbf{W}_1, \mathbf{W}_2)$  is Lipschitz continuous w.r.t.  $\mathbf{W}_1$  and  $\mathbf{W}_2$ . Note that

$$\begin{aligned} & |\ell(\mathbf{W}_1, \mathbf{W}_2, z_i) - \ell(\mathbf{W}'_1, \mathbf{W}_2, z_i)|_2 \\ & \leq \sqrt{2} \left\| \sum_{j=1}^n \left[ g(\tilde{\mathbf{A}}) \right]_{ij} \sigma(\sigma(\mathbf{X}_{j*} \mathbf{W}_1) \mathbf{W}_2) - \sum_{j=1}^n \left[ g(\tilde{\mathbf{A}}) \right]_{ij} \sigma(\sigma(\mathbf{X}_{j*} \mathbf{W}'_1) \mathbf{W}_2) \right\|_2 \\ & \leq \sqrt{2} \sum_{j=1}^n \left| \left[ g(\tilde{\mathbf{A}}) \right]_{ij} \right| \left\| (\sigma(\mathbf{X}_{j*} \mathbf{W}_1) - \sigma(\mathbf{X}_{j*} \mathbf{W}'_1)) \mathbf{W}_2 \right\|_2 \\ & \leq \sqrt{2} \sum_{j=1}^n \left| \left[ g(\tilde{\mathbf{A}}) \right]_{ij} \right| \|\sigma(\mathbf{X}_{j*} \mathbf{W}_1) - \sigma(\mathbf{X}_{j*} \mathbf{W}'_1)\|_2 \|\mathbf{W}_2\| \\ & \leq \sqrt{2} \sum_{j=1}^n \left| \left[ g(\tilde{\mathbf{A}}) \right]_{ij} \right| \|\mathbf{X}_{j*}(\mathbf{W}_1 - \mathbf{W}'_1)\|_2 \|\mathbf{W}_2\| \\ & \leq c_X c_W \sqrt{2} \left\| g(\tilde{\mathbf{A}}) \right\|_\infty \|\text{vec}(\mathbf{W}_1) - \text{vec}(\mathbf{W}'_1)\|. \end{aligned}$$

Besides,

$$\begin{aligned} & |\ell(\mathbf{W}_1, \mathbf{W}_2, z_i) - \ell(\mathbf{W}_1, \mathbf{W}'_2, z_i)|_2 \\ & \leq \sqrt{2} \left\| \sum_{j=1}^n \left[ g(\tilde{\mathbf{A}}) \right]_{ij} (\sigma(\sigma(\mathbf{X}_{j*} \mathbf{W}_1) \mathbf{W}_2) - \sigma(\sigma(\mathbf{X}_{j*} \mathbf{W}_1) \mathbf{W}'_2)) \right\|_2 \\ & \leq \sqrt{2} \sum_{j=1}^n \left| \left[ g(\tilde{\mathbf{A}}) \right]_{ij} \right| \|\sigma(\mathbf{X}_{j*} \mathbf{W}_1)(\mathbf{W}_2 - \mathbf{W}'_2)\|_2 \leq c_X c_W \sqrt{2} \left\| g(\tilde{\mathbf{A}}) \right\|_\infty \|\mathbf{W}_2 - \mathbf{W}'_2\|. \end{aligned}$$

By Lemma A.4, we conclude that  $|\ell(\mathbf{w}) - \ell(\mathbf{w}')| \leq L_{\mathcal{F}} \|\mathbf{w} - \mathbf{w}'\|_2$  holds with  $L_{\mathcal{F}} = 2c_X c_W \left\| g(\tilde{\mathbf{A}}) \right\|_\infty$ . The gradients of  $\ell$  w.r.t.  $\mathbf{W}_1$  and  $\mathbf{W}_2$  are

$$\begin{aligned} \frac{\partial \ell(\mathbf{W}_1, \mathbf{W}_2; z_i)}{\partial \text{vec}[\mathbf{W}_2]} &= \sum_{j=1}^n \left[ g(\tilde{\mathbf{A}}) \right]_{ij} (\sigma'(\mathbf{H}^{(1)} \mathbf{W}_2)_{j*} \odot (\hat{\mathbf{y}}_i - \mathbf{y}_i)) \otimes (\mathbf{X} \mathbf{W}_1)_{j*}, \\ \frac{\partial \ell(\mathbf{W}_1, \mathbf{W}_2; z_i)}{\partial \text{vec}[\mathbf{W}_1]} &= \sum_{j=1}^n \left[ g(\tilde{\mathbf{A}}) \right]_{ij} (\sigma'(\mathbf{X} \mathbf{W}_1)_{j*} \odot (((\hat{\mathbf{y}}_i - \mathbf{y}_i) \odot \sigma'(\mathbf{H}^{(1)} \mathbf{W}_2)_{j*}) \mathbf{W}_2^\top)) \otimes \mathbf{X}_{j*}, \end{aligned}$$

where  $\mathbf{H}^{(1)} = \sigma(\mathbf{X} \mathbf{W}_1)$ . Note that

$$\|\hat{\mathbf{y}}_i(\mathbf{W}_2) - \hat{\mathbf{y}}_i(\mathbf{W}'_2)\|_2 \leq \sum_{j=1}^n \left| \left[ g(\tilde{\mathbf{A}}) \right]_{ij} \right| \|\sigma(\mathbf{H}^{(1)}(\mathbf{W}_2 - \mathbf{W}'_2))\|_2 \leq c_X c_W \left\| g(\tilde{\mathbf{A}}) \right\|_\infty \|\text{vec}[\mathbf{W}_2] - \text{vec}[\mathbf{W}'_2]\|_2.$$

Similarly,

$$\|\hat{\mathbf{y}}_i(\mathbf{W}_1) - \hat{\mathbf{y}}_i(\mathbf{W}'_1)\|_2 \leq c_X c_W \left\| g(\tilde{\mathbf{A}}) \right\|_\infty \|\text{vec}[\mathbf{W}_1] - \text{vec}[\mathbf{W}'_1]\|_2.$$

**Part A.** First we have

$$\begin{aligned} & \left\| \frac{\partial \ell(\mathbf{W}_1, \mathbf{W}_2; z_i)}{\partial \text{vec}[\mathbf{W}_2]}(\mathbf{W}_2) - \frac{\partial \ell(\mathbf{W}_1, \mathbf{W}_2; z_i)}{\partial \text{vec}[\mathbf{W}_2]}(\mathbf{W}'_2) \right\|_2 \\ & \leq \sum_{j=1}^n \left| \left[ g(\tilde{\mathbf{A}}) \right]_{ij} \right| \sqrt{2} \|\mathbf{X}_{j*} \mathbf{W}_1\|_2 \|\sigma'(\mathbf{H}^{(1)} \mathbf{W}_2)_{j*} - \sigma'(\mathbf{H}^{(1)} \mathbf{W}'_2)_{j*}\|_2 \\ & \quad + \sum_{j=1}^n \left| \left[ g(\tilde{\mathbf{A}}) \right]_{ij} \right| \|\mathbf{X}_{j*} \mathbf{W}_1\|_2 \|\hat{\mathbf{y}}_i(\mathbf{W}_2) - \hat{\mathbf{y}}_i(\mathbf{W}'_2)\|_2 \\ & \leq \left\| g(\tilde{\mathbf{A}}) \right\|_\infty \sqrt{|\mathcal{Y}|} c_X^{1+\tilde{\alpha}} c_W^{1+\tilde{\alpha}} \|\text{vec}[\mathbf{W}_2] - \text{vec}[\mathbf{W}'_2]\|_2^{\tilde{\alpha}} + c_X^2 c_W^2 \left\| g(\tilde{\mathbf{A}}) \right\|_\infty^2 \|\text{vec}[\mathbf{W}_2] - \text{vec}[\mathbf{W}'_2]\|_2. \end{aligned}$$

Also,

$$\begin{aligned} & \left\| \frac{\partial \ell(\mathbf{W}_1, \mathbf{W}_2; z_i)}{\partial \text{vec}[\mathbf{W}_2]}(\mathbf{W}_1) - \frac{\partial \ell(\mathbf{W}_1, \mathbf{W}_2; z_i)}{\partial \text{vec}[\mathbf{W}_2]}(\mathbf{W}'_1) \right\|_2 \\ & \leq \sqrt{2} \sum_{j=1}^n \left| \left[ g(\tilde{\mathbf{A}}) \right]_{ij} \right| \|\mathbf{X}_{j*} \mathbf{W}_1\|_2 \|\sigma'(\mathbf{H}^{(1)}(\mathbf{W}_1) \mathbf{W}_2)_{j*} - \sigma'(\mathbf{H}^{(1)}(\mathbf{W}'_1) \mathbf{W}_2)_{j*}\|_2 \\ & \quad + \sqrt{2} \sum_{j=1}^n \left| \left[ g(\tilde{\mathbf{A}}) \right]_{ij} \right| \|\mathbf{X}_{j*}(\mathbf{W}_1 - \mathbf{W}'_1)\|_2 + \sum_{j=1}^n \left| \left[ g(\tilde{\mathbf{A}}) \right]_{ij} \right| \|\hat{\mathbf{y}}_i(\mathbf{W}_1) - \hat{\mathbf{y}}_i(\mathbf{W}'_1)\|_2 \|\mathbf{X}_{j*} \mathbf{W}_1\|_2 \\ & \leq \sqrt{2} P \left\| g(\tilde{\mathbf{A}}) \right\|_\infty c_X^{1+\tilde{\alpha}} c_W^{1+\tilde{\alpha}} \|\text{vec}[\mathbf{W}_1] - \text{vec}[\mathbf{W}'_1]\|_2^{\tilde{\alpha}} \\ & \quad + \left( \sqrt{2} c_X \left\| g(\tilde{\mathbf{A}}) \right\|_\infty + c_X^2 c_W^2 \left\| g(\tilde{\mathbf{A}}) \right\|_\infty^2 \right) \|\text{vec}[\mathbf{W}_1] - \text{vec}[\mathbf{W}'_1]\|_2. \end{aligned}$$

Denote by

$$\begin{aligned} P_{21} &= \left( \sqrt{2} c_X \left\| g(\tilde{\mathbf{A}}) \right\|_\infty + c_X^2 c_W^2 \left\| g(\tilde{\mathbf{A}}) \right\|_\infty^2 \right), P_{22} = c_X^2 c_W^2 \left\| g(\tilde{\mathbf{A}}) \right\|_\infty^2, \\ \tilde{P}_{21} &= \left\| g(\tilde{\mathbf{A}}) \right\|_\infty \sqrt{|\mathcal{Y}|} c_X^{1+\tilde{\alpha}} c_W^{1+\tilde{\alpha}}, \tilde{P}_{22} = \sqrt{2} P \left\| g(\tilde{\mathbf{A}}) \right\|_\infty c_X^{1+\tilde{\alpha}} c_W^{1+\tilde{\alpha}}, \end{aligned}$$

we obtain that

$$\left\| \frac{\partial \ell(\mathbf{w}; z_i)}{\partial \text{vec}[\mathbf{W}_2]} - \frac{\partial \ell(\mathbf{w}'; z_i)}{\partial \text{vec}[\mathbf{W}_2]} \right\|_2 \leq \sum_{i=1}^2 P_{2i} \|\text{vec}[\mathbf{W}_i] - \text{vec}[\mathbf{W}'_i]\|_2 + \tilde{P}_{2i} \|\text{vec}[\mathbf{W}_i] - \text{vec}[\mathbf{W}'_i]\|_2^{\tilde{\alpha}}.$$

**Part B.** We have

$$\begin{aligned} & \left\| \frac{\partial \ell(\mathbf{W}_1, \mathbf{W}_2; z_i)}{\partial \text{vec}[\mathbf{W}_1]}(\mathbf{W}_2) - \frac{\partial \ell(\mathbf{W}_1, \mathbf{W}_2; z_i)}{\partial \text{vec}[\mathbf{W}_1]}(\mathbf{W}'_2) \right\|_2 \\ & \leq c_X c_W \sqrt{2} \sum_{j=1}^n \left| \left[ g(\tilde{\mathbf{A}}) \right]_{ij} \right| \|\sigma'(\mathbf{H}^{(1)} \mathbf{W}_2)_{j*} - \sigma'(\mathbf{H}^{(1)} \mathbf{W}'_2)_{j*}\|_2 + c_X \sqrt{2} \sum_{j=1}^n \left| \left[ g(\tilde{\mathbf{A}}) \right]_{ij} \right| \|\mathbf{W}_2 - \mathbf{W}'_2\|_2 \\ & \quad + c_W \sum_{j=1}^n \left| \left[ g(\tilde{\mathbf{A}}) \right]_{ij} \right| \|\mathbf{X}_{j*}\|_2 \|\hat{\mathbf{y}}_i(\mathbf{W}_2) - \hat{\mathbf{y}}_i(\mathbf{W}'_2)\|_2 \\ & \leq c_X^{1+\tilde{\alpha}} c_W^{1+\tilde{\alpha}} P \sqrt{2} \left\| g(\tilde{\mathbf{A}}) \right\|_\infty \|\text{vec}[\mathbf{W}_2] - \text{vec}[\mathbf{W}'_2]\|_2^{\tilde{\alpha}} \\ & \quad + \left( c_X \sqrt{2} \left\| g(\tilde{\mathbf{A}}) \right\|_\infty + c_X^2 c_W^2 \left\| g(\tilde{\mathbf{A}}) \right\|_\infty^2 \right) \|\text{vec}[\mathbf{W}_2] - \text{vec}[\mathbf{W}'_2]\|_2. \end{aligned}$$

Also, one can find that

$$\begin{aligned}
 & \left\| \frac{\partial \ell(\mathbf{W}_1, \mathbf{W}_2; z_i)}{\partial \text{vec}[\mathbf{W}_1]}(\mathbf{W}_1) - \frac{\partial \ell(\mathbf{W}_1, \mathbf{W}_2; z_i)}{\partial \text{vec}[\mathbf{W}_1]}(\mathbf{W}'_1) \right\|_2 \\
 & \leq c_X c_W \sqrt{2} \sum_{j=1}^n \left[ g(\tilde{\mathbf{A}}) \right]_{ij} \left\| \sigma'(\mathbf{X} \mathbf{W}_1)_{j*} - \sigma'(\mathbf{X} \mathbf{W}'_1)_{j*} \right\|_2 + c_W \sum_{j=1}^n \left[ g(\tilde{\mathbf{A}}) \right]_{ij} \|\mathbf{X}_{j*}\|_2 \|\hat{\mathbf{y}}_i(\mathbf{W}_1) - \hat{\mathbf{y}}_i(\mathbf{W}'_1)\|_2 \\
 & \quad + c_X c_W \sqrt{2} \sum_{j=1}^n \left[ g(\tilde{\mathbf{A}}) \right]_{ij} \left\| \sigma'(\mathbf{H}^{(1)}(\mathbf{W}_1) \mathbf{W}_2) - \sigma'(\mathbf{H}^{(1)}(\mathbf{W}'_1) \mathbf{W}_2) \right\|_2 \\
 & \leq c_X c_W P \sqrt{2} \sum_{j=1}^n \left[ g(\tilde{\mathbf{A}}) \right]_{ij} \|\mathbf{X}_{j*}(\mathbf{W}_1 - \mathbf{W}'_1)\|_2^{\tilde{\alpha}} + c_X c_W P \sqrt{2} \sum_{j=1}^n \left[ g(\tilde{\mathbf{A}}) \right]_{ij} \left\| (\mathbf{H}^{(1)}(\mathbf{W}_1) - \mathbf{H}^{(1)}(\mathbf{W}'_1)) \mathbf{W}_2 \right\|_2^{\tilde{\alpha}} \\
 & \leq [c_X^{1+\tilde{\alpha}} c_W + c_X^{1+\tilde{\alpha}} c_W^{1+\tilde{\alpha}}] P \sqrt{2} \left\| g(\tilde{\mathbf{A}}) \right\|_{\infty} \|\text{vec}[\mathbf{W}_1] - \text{vec}[\mathbf{W}'_1]\|_2^{\tilde{\alpha}} + c_X^2 c_W^2 \left\| g(\tilde{\mathbf{A}}) \right\|_{\infty}^2 \|\text{vec}[\mathbf{W}_1] - \text{vec}[\mathbf{W}'_1]\|_2.
 \end{aligned}$$

Denote by

$$\begin{aligned}
 P_{11} &= c_X^2 c_W^2 \left\| g(\tilde{\mathbf{A}}) \right\|_{\infty}^2, \quad P_{12} = \left( c_X \sqrt{2} \left\| g(\tilde{\mathbf{A}}) \right\|_{\infty} + c_X^2 c_W^2 \left\| g(\tilde{\mathbf{A}}) \right\|_{\infty}^2 \right), \\
 \tilde{P}_{11} &= c_X^{1+\tilde{\alpha}} c_W^{1+\tilde{\alpha}} P \sqrt{2} \left\| g(\tilde{\mathbf{A}}) \right\|_{\infty}, \quad \tilde{P}_{12} = [c_X^{1+\tilde{\alpha}} c_W + c_X^{1+\tilde{\alpha}} c_W^{1+\tilde{\alpha}}] P \sqrt{2} \left\| g(\tilde{\mathbf{A}}) \right\|_{\infty},
 \end{aligned}$$

we obtain that

$$\left\| \frac{\partial \ell(\mathbf{w}; z_i)}{\partial \text{vec}[\mathbf{W}_1]} - \frac{\partial \ell(\mathbf{w}'; z_i)}{\partial \text{vec}[\mathbf{W}_1]} \right\|_2 \leq \sum_{i=1}^2 P_{1i} \|\text{vec}[\mathbf{W}_i] - \text{vec}[\mathbf{W}'_i]\|_2 + \sum_{i=1}^2 \tilde{P}_{1i} \|\text{vec}[\mathbf{W}_i] - \text{vec}[\mathbf{W}'_i]\|_2^{\tilde{\alpha}}.$$

By Lemma A.4, we conclude that  $\|\nabla \ell(\mathbf{w}) - \nabla \ell(\mathbf{w}')\|_2 \leq P_{\mathcal{F}} \max\{\|\mathbf{w} - \mathbf{w}'\|_2, \|\mathbf{w} - \mathbf{w}'\|_2^{\tilde{\alpha}}\}$  holds where  $\mathbf{w} = [\text{vec}[\mathbf{W}_1]; \text{vec}[\mathbf{W}_2]]$ .

### B.2.5. PROOF OF PROPOSITION 4.15

We first show that the objective  $\ell(\mathbf{W}_1, \mathbf{W}_2, \gamma)$  is Lipschitz continuous w.r.t.  $\mathbf{W}_1, \mathbf{W}_2$  and  $\gamma$ . Note that

$$\begin{aligned}
 & |\ell(\mathbf{W}_1, \mathbf{W}_2, \gamma) - \ell(\mathbf{W}_1, \mathbf{W}_2, \gamma')| \\
 & \leq \sqrt{2} \left\| \sum_{j=1}^n \left[ g(\tilde{\mathbf{A}}, \gamma) \right]_{ij} \sigma(\sigma(\mathbf{X}_{j*} \mathbf{W}_1) \mathbf{W}_2) - \sum_{j=1}^n \left[ g(\tilde{\mathbf{A}}, \gamma') \right]_{ij} \sigma(\sigma(\mathbf{X}_{j*} \mathbf{W}_1) \mathbf{W}_2) \right\|_2 \\
 & = \sqrt{2} \left\| \sum_{k=0}^K (\gamma_k - \gamma'_k) \left( \sum_{j=1}^n \left[ \tilde{\mathbf{A}}^k \right]_{ij} \sigma(\sigma(\mathbf{X}_{j*} \mathbf{W}_1) \mathbf{W}_2) \right) \right\|_2 \\
 & \leq \sqrt{2} \|\gamma - \gamma'\|_2 \sqrt{\sum_{k=0}^K \left\| \sum_{j=1}^n \left[ \tilde{\mathbf{A}}^k \right]_{ij} \sigma(\sigma(\mathbf{X}_{j*} \mathbf{W}_1) \mathbf{W}_2) \right\|_2^2} \\
 & \leq \sqrt{2} \|\gamma - \gamma'\|_2 \sum_{k=0}^K \left\| \sum_{j=1}^n \left[ \tilde{\mathbf{A}}^k \right]_{ij} \sigma(\sigma(\mathbf{X}_{j*} \mathbf{W}_1) \mathbf{W}_2) \right\|_2 \leq \sqrt{2} c_X c_W^2 \left( \sum_{k=0}^K \left\| \tilde{\mathbf{A}}^k \right\|_{\infty} \right) \|\gamma - \gamma'\|_2.
 \end{aligned}$$

Denote by

$$L_{\mathcal{F}} = \sqrt{2c_X^2 c_W^4 \left( \sum_{k=0}^K \left\| \tilde{\mathbf{A}}^k \right\|_{\infty} \right)^2 + 4c_X^2 c_W^2 \left\| g(\tilde{\mathbf{A}}, \gamma) \right\|_{\infty}^2},$$

we conclude that  $|\ell(\mathbf{w}) - \ell(\mathbf{w}')| \leq L_{\mathcal{F}} \|\mathbf{w} - \mathbf{w}'\|_2$  holds. Then we discuss the smoothness of this model. The gradients of  $\ell(\mathbf{W}_1, \mathbf{W}_2, \gamma)$  w.r.t.  $\mathbf{W}_1$ ,  $\mathbf{W}_2$ , and  $\gamma$  are

$$\begin{aligned} \frac{\partial \ell(\mathbf{W}_1, \mathbf{W}_2, \gamma; z_i)}{\partial \text{vec}[\mathbf{W}_2]} &= \sum_{j=1}^n \left[ g(\tilde{\mathbf{A}}, \gamma) \right]_{ij} (\sigma'(\sigma(\mathbf{X}\mathbf{W}_1)\mathbf{W}_2)_{j*} \odot (\hat{\mathbf{y}}_i - \mathbf{y}_i)) \otimes (\mathbf{X}\mathbf{W}_1)_{j*}, \\ \frac{\partial \ell(\mathbf{W}_1, \mathbf{W}_2, \gamma; z_i)}{\partial \text{vec}[\mathbf{W}_1]} &= \sum_{j=1}^n \left[ g(\tilde{\mathbf{A}}, \gamma) \right]_{ij} (\sigma'(\mathbf{X}\mathbf{W}_1)_{j,:} \odot (((\hat{\mathbf{y}}_i - \mathbf{y}_i) \odot \sigma'(\sigma(\mathbf{X}\mathbf{W}_1)\mathbf{W}_2)_{j*}) \mathbf{W}_2^\top)) \otimes \mathbf{X}_{j*}, \\ \frac{\partial \ell(\mathbf{W}_1, \mathbf{W}_2, \gamma; z_i)}{\partial \gamma} &= (\hat{\mathbf{y}}_i - \mathbf{y}_i) \left[ \sum_{j=1}^n [\tilde{\mathbf{A}}^0]_{ij} \mathbf{H}_{j*}^\top, \dots, \sum_{j=1}^n [\tilde{\mathbf{A}}^K]_{ij} \mathbf{H}_{j*}^\top \right], \end{aligned}$$

where  $\mathbf{H} = \sigma(\sigma(\mathbf{X}\mathbf{W}_1)\mathbf{W}_2)$ . Note that

$$\|\hat{\mathbf{y}}_i(\gamma) - \hat{\mathbf{y}}_i(\gamma')\|_2 \leq c_X c_W^2 \left( \sum_{k=0}^K \|\tilde{\mathbf{A}}^k\|_\infty \right) \|\gamma - \gamma'\|_2.$$

Then one can find that

$$\begin{aligned} & \left\| \frac{\partial \ell(\mathbf{W}_1, \mathbf{W}_2, \gamma; z_i)}{\partial \text{vec}[\mathbf{W}_2]}(\gamma) - \frac{\partial \ell(\mathbf{W}_1, \mathbf{W}_2, \gamma; z_i)}{\partial \text{vec}[\mathbf{W}_2]}(\gamma') \right\|_2 \\ & \leq \left\| \sum_{k=0}^K (\gamma_k - \gamma'_k) \left( \sum_{j=1}^n [\tilde{\mathbf{A}}^k]_{ij} (\sigma'(\sigma(\mathbf{X}\mathbf{W}_1)\mathbf{W}_2)_{j*} \odot (\hat{\mathbf{y}}_i - \mathbf{y}_i)) \otimes (\mathbf{X}\mathbf{W}_1)_{j*} \right) \right\| \\ & \quad + \sum_{j=1}^n \left| [g(\tilde{\mathbf{A}}, \gamma)]_{ij} \right| \|\hat{\mathbf{y}}_i(\gamma) - \hat{\mathbf{y}}_i(\gamma')\|_2 \|(\mathbf{X}\mathbf{W}_1)_{j*}\|_2 \\ & \leq (\sqrt{2} + c_X c_W \|g(\tilde{\mathbf{A}}, \gamma)\|_\infty) c_X c_W^2 \left( \sum_{k=0}^K \|\tilde{\mathbf{A}}^k\|_\infty \right) \|\gamma - \gamma'\|_2. \end{aligned}$$

Denote by

$$\begin{aligned} P_{21} &= c_X^2 c_W^2 \|g(\tilde{\mathbf{A}}, \gamma)\|_\infty^2, \quad P_{22} = \left( \sqrt{2} c_X \|g(\tilde{\mathbf{A}}, \gamma)\|_\infty + c_X^2 c_W^2 \|g(\tilde{\mathbf{A}}, \gamma)\|_\infty^2 \right), \\ P_{23} &= \left( \sqrt{2} + c_X c_W \|g(\tilde{\mathbf{A}}, \gamma)\|_\infty \right) c_X c_W^2 \left( \sum_{k=0}^K \|\tilde{\mathbf{A}}^k\|_\infty \right), \\ \tilde{P}_{21} &= \|g(\tilde{\mathbf{A}}, \gamma)\|_\infty \sqrt{|\mathcal{Y}|} c_X^{1+\tilde{\alpha}} c_W^{1+\tilde{\alpha}}, \quad \tilde{P}_{22} = \sqrt{2} P \|g(\tilde{\mathbf{A}}, \gamma)\|_\infty c_X^{1+\tilde{\alpha}} c_W^{1+\tilde{\alpha}}, \quad \tilde{P}_{23} = 0, \end{aligned}$$

we obtain that

$$\begin{aligned} \left\| \frac{\partial \ell(\mathbf{w}; z_i)}{\partial \text{vec}[\mathbf{W}_1]} - \frac{\partial \ell(\mathbf{w}'; z_i)}{\partial \text{vec}[\mathbf{W}_1]} \right\|_2 & \leq \sum_{i=1}^2 P_{2i} \|\text{vec}[\mathbf{W}_i] - \text{vec}[\mathbf{W}'_i]\|_2 + \sum_{i=1}^2 \tilde{P}_{2i} \|\text{vec}[\mathbf{W}_i] - \text{vec}[\mathbf{W}'_i]\|_2^{\tilde{\alpha}} \\ & \quad + P_{23} \|\gamma - \gamma'\|_2 + \tilde{P}_{23} \|\gamma - \gamma'\|_2^{\tilde{\alpha}}. \end{aligned}$$

Besides,

$$\begin{aligned} & \left\| \frac{\partial \ell(\mathbf{W}_1, \mathbf{W}_2, \gamma; z_i)}{\partial \text{vec}[\mathbf{W}_1]}(\gamma) - \frac{\partial \ell(\mathbf{W}_1, \mathbf{W}_2, \gamma; z_i)}{\partial \text{vec}[\mathbf{W}_1]}(\gamma') \right\| \\ & = \left\| \sum_{k=0}^K (\gamma_k - \gamma'_k) \left( \sum_{j=1}^n [\tilde{\mathbf{A}}^k]_{ij} (\sigma'(\mathbf{X}\mathbf{W}_1)_{j,:} \odot (((\hat{\mathbf{y}}_i - \mathbf{y}_i) \odot \sigma'(\sigma(\mathbf{X}\mathbf{W}_1)\mathbf{W}_2)_{j*}) \mathbf{W}_2^\top)) \otimes \mathbf{X}_{j*} \right) \right\|_2 \\ & \quad + \sum_{j=1}^n c_X c_W \left| [g(\tilde{\mathbf{A}}, \gamma)]_{ij} \right| \|\hat{\mathbf{y}}_i(\gamma) - \hat{\mathbf{y}}_i(\gamma')\|_2 \\ & \leq (\sqrt{2} + c_X c_W \|g(\tilde{\mathbf{A}}, \gamma)\|_\infty) c_X c_W^2 \left( \sum_{k=0}^K \|\tilde{\mathbf{A}}^k\|_\infty \right) \|\gamma - \gamma'\|_2. \end{aligned}$$

Denote by

$$\begin{aligned}
 P_{11} &= c_X^2 c_W^2 \left\| g(\tilde{\mathbf{A}}, \gamma) \right\|_\infty^2, \quad P_{12} = \left( c_X \sqrt{2} \left\| g(\tilde{\mathbf{A}}, \gamma) \right\|_\infty + c_X^2 c_W^2 \left\| g(\tilde{\mathbf{A}}, \gamma) \right\|_\infty^2 \right), \\
 P_{13} &= \left( \sqrt{2} + c_X c_W \left\| g(\tilde{\mathbf{A}}, \gamma) \right\|_\infty \right) c_X c_W^2 \left( \sum_{k=0}^K \left\| \tilde{\mathbf{A}}^k \right\|_\infty \right) \\
 \tilde{P}_{11} &= \sqrt{2} P \left[ c_X^{1+\tilde{\alpha}} c_W + c_X^{1+\tilde{\alpha}} c_W^{1+\tilde{\alpha}} \right], \quad \tilde{P}_{12} = c_X^{1+\tilde{\alpha}} c_W^{1+\tilde{\alpha}} P \sqrt{2} \left\| g(\tilde{\mathbf{A}}, \gamma) \right\|_\infty, \quad \tilde{P}_{13} = 0,
 \end{aligned}$$

we obtain that

$$\begin{aligned}
 \left\| \frac{\partial \ell(\mathbf{w}; z_i)}{\partial \text{vec}[\mathbf{W}_1]} - \frac{\partial \ell(\mathbf{w}'; z_i)}{\partial \text{vec}[\mathbf{W}_1]} \right\|_2 &\leq \sum_{i=1}^2 P_{1i} \left\| \text{vec}[\mathbf{W}_i] - \text{vec}[\mathbf{W}'_i] \right\|_2 + \sum_{i=1}^2 \tilde{P}_{1i} \left\| \text{vec}[\mathbf{W}_i] - \text{vec}[\mathbf{W}'_i] \right\|_2^{\tilde{\alpha}} \\
 &\quad + P_{13} \left\| \gamma - \gamma' \right\|_2 + \tilde{P}_{13} \left\| \gamma - \gamma' \right\|_2^{\tilde{\alpha}}.
 \end{aligned}$$

Lastly, since

$$\begin{aligned}
 &\left\| \frac{\partial \ell(\mathbf{W}_1, \mathbf{W}_2, \gamma; z_i)}{\partial \gamma}(\mathbf{W}_2) - \frac{\partial \ell(\mathbf{W}_1, \mathbf{W}_2, \gamma; z_i)}{\partial \gamma}(\mathbf{W}'_2) \right\|_2 \\
 &\leq \left\| \hat{\mathbf{y}}_i - \mathbf{y}_i \right\|_2 \sqrt{\sum_{k=0}^K \left\| \sum_{j=1}^n [\tilde{\mathbf{A}}^k]_{ij} (\sigma'(\sigma(\mathbf{X}_j \mathbf{W}_1) \mathbf{W}_2) - \sigma'(\sigma(\mathbf{X}_j \mathbf{W}_1) \mathbf{W}'_2)) \right\|_2^2} \\
 &\quad + \left\| \hat{\mathbf{y}}_i(\mathbf{W}_2) - \hat{\mathbf{y}}_i(\mathbf{W}'_2) \right\|_2 \sqrt{\sum_{k=0}^K \left\| \sum_{j=1}^n [\tilde{\mathbf{A}}^k]_{ij} \sigma'(\sigma(\mathbf{X}_j \mathbf{W}_1) \mathbf{W}_2) \right\|_2^2} \\
 &\leq \sqrt{2} P c_X^{\tilde{\alpha}} c_W^{\tilde{\alpha}} \left( \sum_{k=0}^K \left\| \tilde{\mathbf{A}}^k \right\|_\infty \right) \left\| \text{vec}[\mathbf{W}_2] - \text{vec}[\mathbf{W}'_2] \right\|_2^{\tilde{\alpha}} \\
 &\quad + c_X^2 c_W^3 \left\| g(\tilde{\mathbf{A}}, \gamma) \right\|_\infty \left( \sum_{k=0}^K \left\| \tilde{\mathbf{A}}^k \right\|_\infty \right) \left\| \text{vec}[\mathbf{W}_2] - \text{vec}[\mathbf{W}'_2] \right\|_2.
 \end{aligned}$$

Similarly, we have

$$\begin{aligned}
 &\left\| \frac{\partial \ell(\mathbf{W}_1, \mathbf{W}_2, \gamma; z_i)}{\partial \gamma}(\mathbf{W}_1) - \frac{\partial \ell(\mathbf{W}_1, \mathbf{W}_2, \gamma; z_i)}{\partial \gamma}(\mathbf{W}'_1) \right\|_2 \\
 &\leq \sqrt{2} P c_X^{\tilde{\alpha}} c_W^{\tilde{\alpha}} \left( \sum_{k=0}^K \left\| \tilde{\mathbf{A}}^k \right\|_\infty \right) \left\| \text{vec}[\mathbf{W}_1] - \text{vec}[\mathbf{W}'_1] \right\|_2^{\tilde{\alpha}} \\
 &\quad + c_X^2 c_W^3 \left\| g(\tilde{\mathbf{A}}, \gamma) \right\|_\infty \left( \sum_{k=0}^K \left\| \tilde{\mathbf{A}}^k \right\|_\infty \right) \left\| \text{vec}[\mathbf{W}_1] - \text{vec}[\mathbf{W}'_1] \right\|_2,
 \end{aligned}$$

and

$$\begin{aligned}
 &\left\| \frac{\partial \ell(\mathbf{W}_1, \mathbf{W}_2, \gamma; z_i)}{\partial \gamma}(\gamma) - \frac{\partial \ell(\mathbf{W}_1, \mathbf{W}_2, \gamma; z_i)}{\partial \gamma}(\gamma') \right\|_2 \\
 &\leq \left\| \hat{\mathbf{y}}_i(\gamma) - \hat{\mathbf{y}}_i(\gamma') \right\|_2 \sqrt{\sum_{k=0}^K \left\| \sum_{j=1}^n [\tilde{\mathbf{A}}^k]_{ij} \sigma(\sigma(\mathbf{X}_j \mathbf{W}_1) \mathbf{W}_2) \right\|_2^2} \leq c_X^2 c_W^4 \left( \sum_{k=0}^K \left\| \tilde{\mathbf{A}}^k \right\|_\infty \right)^2 \left\| \gamma - \gamma' \right\|_2.
 \end{aligned}$$

Denote by

$$\begin{aligned}
 P_{31} &= P_{32} = c_X^2 c_W^3 \left\| g(\tilde{\mathbf{A}}, \gamma) \right\|_\infty \left( \sum_{k=0}^K \left\| \tilde{\mathbf{A}}^k \right\|_\infty \right), \quad P_{33} = c_X^2 c_W^4 \left( \sum_{k=0}^K \left\| \tilde{\mathbf{A}}^k \right\|_\infty \right)^2, \\
 \tilde{P}_{31} &= \tilde{P}_{32} = \sqrt{2} P c_X^{\tilde{\alpha}} c_W^{\tilde{\alpha}} \left( \sum_{k=0}^K \left\| \tilde{\mathbf{A}}^k \right\|_\infty \right), \quad \tilde{P}_{33} = 0,
 \end{aligned}$$

we obtain that

$$\left\| \frac{\partial \ell(\mathbf{w}; z_i)}{\partial \gamma} - \frac{\partial \ell(\mathbf{w}'; z_i)}{\partial \gamma} \right\|_2 \leq \sum_{i=1}^2 P_{3i} \|\text{vec}[\mathbf{W}_i] - \text{vec}[\mathbf{W}'_i]\|_2 + \sum_{i=1}^2 \tilde{P}_{3i} \|\text{vec}[\mathbf{W}_i] - \text{vec}[\mathbf{W}'_i]\|_2^{\tilde{\alpha}}$$

$$+ P_{33} \|\gamma - \gamma'\|_2 + \tilde{P}_{33} \|\gamma - \gamma'\|_2^{\tilde{\alpha}}.$$

By Lemma A.4, we conclude that  $\|\nabla \ell(\mathbf{w}) - \nabla \ell(\mathbf{w}')\|_2 \leq P_{\mathcal{F}} \max\{\|\mathbf{w} - \mathbf{w}'\|_2, \|\mathbf{w} - \mathbf{w}'\|_2^{\tilde{\alpha}}\}$  holds where  $\mathbf{w} = [\text{vec}[\mathbf{W}_1]; \text{vec}[\mathbf{W}_2]; \gamma]$ .

## C. Experiments Details

For GCN, GAT, SGC, APPNP and GCNII, we adopt the official PyTorch Geometric library implementations (Fey & Lenssen, 2019). For GPR-GNN, we adopt the released codes <sup>1</sup> with commit number 2507f10. Following (Cong et al., 2021), we remove all dropout layers and adopt the Adam optimizer with default setting. The batch size is set to 512 and the number of hidden units are set to 64 for all baseline models.  $K$  is set to 10 for APPNP and GPR-GNN.

---

<sup>1</sup><https://github.com/jianhao2016/GPRGNN>