

# Theoretical Learning Performance of Graph Networks: the Impact of Jumping Connections and Layer-wise Sparsification

Anonymous authors

Paper under double-blind review

## Abstract

Jumping connections enable Graph Convolutional Networks (GCNs) to overcome over-smoothing, while graph sparsification reduces computational demands by selecting a sub-matrix of the graph adjacency matrix during neighborhood aggregation. Learning GCNs with graph sparsification has shown empirical success across various applications, but a theoretical understanding of the generalization guarantees remains limited, with existing analyses ignoring either graph sparsification or jumping connections. This paper presents the first learning dynamics and generalization analysis of GCNs with jumping connections using graph sparsification. Our analysis demonstrates that the generalization accuracy of the learned model closely approximates the highest achievable accuracy within a broad class of target functions dependent on the proposed sparse effective adjacency matrix  $A^*$ . Thus, graph sparsification maintains generalization performance when  $A^*$  accurately models data correlations. We reveal that jumping connections lead to different sparsification requirements across layers. In a two-hidden-layer GCN, the generalization is more affected by the sparsified matrix deviations from  $A^*$  of the first layer than the second layer. To the best of our knowledge, this marks the first theoretical characterization of jumping connections' role in sparsification requirements. We validate our theoretical results on benchmark datasets in deep GCNs.

## 1 Introduction

Graph neural networks (GNNs) outperform traditional machine learning techniques when learning graph-structured data, that comprises a collection of features linked with nodes and a graph representing the correlation of the features. As one of the most popular variants of GNN, Graph Convolutional Networks (GCNs) (Kipf & Welling, 2017) perform the convolution operations on graphs by aggregating neighboring nodes to update the feature presentation of every node. GCNs have demonstrated great empirical success such as text classification (Norcliffe-Brown et al., 2018; Zhang et al., 2018) and recommendation systems (Wu et al., 2019; Ying et al., 2018). Because GCNs are easy and computationally efficient to implement, they are the preferred choice for large-scale graph training Duan et al. (2022); Zhang et al. (2022).

As the depth of vanilla GCNs increases, there is a tendency for the node representations to converge toward a common value, a phenomenon known as “over-smoothing” (Li et al., 2018). A widely adopted mitigation approach is to incorporate jumping connections, which allow features to bypass intermediate layers and directly contribute to future layers' output (Li et al., 2019; Xu et al., 2018b). Moreover, jumping connections are shown to accelerate the training process (Xu et al., 2021). Jumping connections have thus become an essential component of GCNs.

Processing large-scale graphs can be computationally demanding, particularly when dealing with the recursive neighborhood integration in GCNs. To alleviate this computational burden, graph sampling or sparsification methods select a subset of nodes or edges from the original graph when computing neighborhood aggregation. Various graph sampling approaches have been developed, including node sampling methods like GraphSAGE (Hamilton et al., 2017), layer-wise sampling like FastGCN (Chen et al., 2018),

subgraph sampling methods like Graphsaint (Zeng et al., 2020). The graph sparsification methods (Li et al., 2020; Chen et al., 2021; You et al., 2020; Liu et al., 2023) usually co-optimize weights and sparsified masks to find optimal sparse graphs and remove the task-irrelevant edges. Ioannidis et al. (2020); Zeng et al. (2020); Srinivasa et al. (2020) introduce simple pruning methods that remove less significant edges without needing complex iterative training.

Although GCNs have demonstrated superior empirical success, their theoretical foundation remains relatively underdeveloped. Some works analyze the expressive power of GCNs in terms of the functions they can represent Morris et al. (2019); Cong et al. (2021); Oono & Suzuki (2019); Xu et al. (2019); Chen et al. (2019), while some works characterize the generalization gap which measures the gap between the training accuracy and test accuracy Esser et al. (2021); Liao et al. (2020); Tang & Liu (2023b). All these works ignore training dynamics, i.e., they do not characterize how to train a model to achieve great expressive power or a small generalization gap. Some works exploit the neural tangent kernel (NTK) technique to analyze the training dynamics of stochastic gradient descent and generalization performance simultaneously Yang et al. (2023). These analyses apply to deep neural networks, only when the network is impractically overparameterized, i.e., the number of neurons is either infinite Du et al. (2019a) or polynomial in the total number of nodes Qin et al. (2023).

Li et al. (2022a); Zhang et al. (2023b); Tang & Liu (2023a) analyze the training dynamics of SGD for GCNs with sparsification and prove that the learned model is guaranteed to achieve desirable generalization. These analyses focus on two-hidden-layer GCNs, but the learning problem is already a nonconvex problem for these shallow networks. However, their network architectures exclude jumping connections. Xu et al. (2021) investigates training dynamics of jumping connections in multi-layer GCNs, but that paper does not provide generalization results and only considers linear activation functions.

To the best of our knowledge, this paper provides the first theoretical analysis of the training dynamics and generalization performance for two-hidden-layer GCNs with jumping-connection using graph sparsification. Our method focuses on the interaction of the jumping-connection and the intermediate layer and explains how jumping-connection influences training and graph sparsification across layers. We consider the semi-supervised node regression problem, where given all node features and partial labels, the objective is to predict the unknown node labels. Our major results include:

- (1) We analyze training two-hidden-layer GCNs by stochastic gradient descent (SGD) with a jumping connection, using a pruning method that prefers the large weight edges from the adjacency matrix  $A$ . Our analysis demonstrates that the generalization accuracy of the learned model approximates the highest achievable accuracy within a broad class of target functions, which map input features to labels. Each target function is a sum of a simpler base function that contributes significantly to the output and a more complicated composite function that has a comparatively smaller impact on the output. This class encompasses a wide range of functions, including two-hidden-layer GCNs with jumping connections.
- (2) This paper extends the concept of the sparse effective adjacency matrix, denoted as  $A^*$ , which is first introduced in Li et al. (2022a) for GCNs lacking jumping connections. This paper shows that  $A^*$  also characterizes the influence of graph sparsification in GCNs with jumping connections. We find that the adjacency matrix  $A$  of a graph often includes redundant information, suggesting that an effectively sparse graph can perform as well as or even surpass  $A$  in training GCNs. Then the goal of graph pruning shifts from minimizing the difference between sampled adjacency matrix, denoted by  $A^s$ , and  $A$  to minimizing the difference between  $A^s$  and  $A^*$ . Consequently, even when the pruned adjacency matrix  $A^s$  is very sparse and significantly deviates from  $A$ , as long as there exists an  $A^*$  closely enough to  $A^s$ , graph sparsification does not compromise the model’s generalization performance.
- (3) This paper theoretically demonstrates that, owing to the presence of jumping connections, sparsifying in different layers has different impacts on the model’s output. Specifically, the first layer connects directly to the output through the jumping connection, and as a result, the deviation of the sampled matrix from the sparse effective matrix  $A^*$  has a more significant effect than the deviation of the second layer. In contrast, the second layer influences the output through a composite function that contributes less significantly, allowing for more substantial deviations from  $A^*$  in the process without compromising error rates. To the best of our knowledge, this is the first theoretical characterization of how jumping connections influence sparsification

requirements, while previous analyses such as (Li et al., 2022a) assume that the sparsification approach remains consistent across different layers. Besides, our experiments on the deep-layer Jumping Knowledge Network GCN, demonstrate the significant impact of graph sampling in shallow layers compared to deeper layers. This empirical evidence supports our theoretical claims and the relevance of our two-layer model analysis in understanding deeper GCN architectures.

### 1.1 Related works

**Other theoretical analysis of GNNs** focus on expressive power and convergence analysis. Xu et al. (2018a); Morris et al. (2019) show the power of 1-hop message passing is upper bounded by 1-WL test. Feng et al. (2022); Wang & Zhang (2022) extend the analysis to k-hop message passing neural networks and spectral GNNs. Zhang et al. (2023a) explores the expressive power of GNNs from the perspective of graph biconnectivity. Oono & Suzuki (2020); Ramezani et al. (2020); Cong et al. (2021) investigates the optimization of GNN training.

**Generalization analyses of Neural Networks (NNs).** Various approaches have been developed to analyze the generalization of feedforward NNs. The neural tangent kernel (NTK) approach shows that overparameterized networks can be approximated by kernel methods in the limiting case (Jacot et al., 2018; Du et al., 2019b). The model estimation approach assumes the existence of a ground-truth one-hidden-layer model with desirable generalization and estimates the model parameters using the training data (Zhong et al., 2017; Zhang et al., 2020; Li et al., 2022b). The feature learning approach analyzes how a shallow NN learns important features during training and thus achieves desirable generalization (Li & Liang, 2018; Allen-Zhu & Li, 2022; 2023). All works ignore the jumping connection except Allen-Zhu & Li (2019), which analyzes the generalization of two-hidden-layer ResNet. Our analysis builds upon Allen-Zhu & Li (2019) and extends to GCNs with graph sparsification.

**Various GNN sparsification methods.** Node-wise sampling (Hamilton et al., 2017) randomly selects nodes and their multi-hop neighbors to create a localized subgraph. Layer-wise sampling (Chen et al., 2018; Zou et al., 2019; Huang et al., 2018) sample a fixed number of nodes in each layer. Subgraph-based sampling (Zeng et al., 2020; Chiang et al., 2019) generates subgraphs by sampling nodes and edges. As for graph sparsification, SGCN (Li et al., 2020) introduces the alternating direction method of multipliers (ADMM) to sparsify the adjacency matrix. UGS (Chen et al., 2021), Early-Bird (You et al., 2020), and ICPG (Sui et al., 2022) design a pruning strategy to sparsify the graph based on the lottery ticket hypothesis. CGP (Liu et al., 2023) proposes a graph gradual pruning framework to reduce the computational cost.

## 2 Training GCNs with Layer-wise Graph Sparsification: Summary of Main Components

### 2.1 GCN Learning Setup

Let  $\mathcal{G} = \{\mathcal{V}, \mathcal{E}\}$  represent an undirected graph, where  $\mathcal{V}$  is the set of nodes with  $|\mathcal{V}| = N$  nodes and  $\mathcal{E}$  is the set of edges.  $\Delta$  is the maximum degree of  $\mathcal{G}$ . An adjacency matrix  $\tilde{A} \in \mathbb{R}^{N \times N}$  is defined to describe the overall graph topology where  $\tilde{A}(i, j) = 1$  if  $(v_i, v_j) \in \mathcal{E}$  else  $\tilde{A}(i, j) = 0$ .  $A$  denotes the normalized adjacency matrix with  $A = D^{-\frac{1}{2}}(\tilde{A} + I)D^{-\frac{1}{2}}$  where  $D$  is the degree matrix with diagonal elements  $D_{i,i} = \sum_j \tilde{A}(i, j)$ . Each element  $A_{ij}$  of the matrix  $A$  represents the normalized weight of the edge between nodes  $i$  and  $j$ .  $a_n$  denotes the  $n$ th column of the  $A$ . Let  $X \in \mathbb{R}^{d \times N}$  denote the feature matrix of  $N$  nodes, where  $\tilde{x}_n \in \mathbb{R}^d$  denotes the feature of node  $n$ . Assume  $\|\tilde{x}_n\| = 1$  for all  $n$  without loss of generality.  $y_n \in \mathbb{R}^k$  represents the label of node  $n$ . Let  $\Omega \subset \mathcal{V}$  denote the set of labeled nodes. Given  $X$  and partial labels in  $\Omega$ , the objective of semi-supervised node-regression is to predict the unknown labels in  $\mathcal{V}/\Omega$ .

We consider training a two-hidden-layer GCN with a single jumping connection, where the function out :  $\mathbb{R}^{d \times N} \times \mathbb{R}^{N \times N} \rightarrow \mathbb{R}^{k \times N}$  with

$$\text{out}(X, A; W, U) = C\sigma(WXA) + C\sigma(U\sigma(WXA)A) \quad (1)$$

where  $\sigma(\cdot)$  applies the ReLU activation  $\text{ReLU}(x) = \max(x, 0)$  to each entry,  $W \in \mathbb{R}^{m \times d}$ ,  $U \in \mathbb{R}^{m \times m}$ , and  $C \in \mathbb{R}^{k \times m}$  denote the first hidden-layer, second hidden-layer, and output layer weights, respectively. We

only train  $W$  and  $U$ . The output of the  $n$ th node can be written as  $\text{out}_n : \mathbb{R}^{d \times N} \times \mathbb{R}^N \rightarrow \mathbb{R}^k$  is

$$\text{out}_n(X, A; W, U) = C\sigma(WXa_n) + C\sigma(U\sigma(WXA)a_n) \quad (2)$$

We focus on the  $\ell_2$  regression task and the prediction loss of the  $n$ th node is defined as

$$\text{Obj}_n(X, A, y_n; W, U) = \frac{1}{2} \|y_n - \text{out}_n(X, A; W, U)\|_2^2 \quad (3)$$

The learning problem solves the following empirical risk minimization problem:

$$\min_{W, U} L_\Omega(W, U) = \frac{1}{|\Omega|} \sum_{n \in \Omega} \text{Obj}_n(X, A, y_n; W, U) \quad (4)$$

## 2.2 Training with stochastic gradient descent and graph sparsification

The recursive neighborhood aggregation through multiplying the feature matrix with  $A$  is costly in both computation and memory. Graph sparsification prunes the graph adjacency matrix  $A$  to reduce the computation and memory requirement. For example, one common theme of various edge sampling or sparsification methods is to retain the large weight edges  $A_{ij}$  from the adjacency matrix  $A$  in  $A^s$  while pruning small weights (Chen et al., 2018; Zeng et al., 2020). To further reduce the computation, layer-wise sampling is also employed that uses different sampling rates in different layers, see, e.g., Chen et al. (2018).

We allow the sparsification methods with different parameter settings in different layers. Specifically, In the  $t$ th iteration, let  $A^{1t}$  and  $A^{2t}$  denote the sparsified adjacency matrix  $A$  in the first and second hidden layers, respectively.

In Algorithm 1, (4) is solved by the stochastic gradient descent (SGD) method starting from random initialization. In each iteration, the gradient of the prediction loss of one randomly selected node is used to approximate the gradient of  $L_\Omega$ . Let  $W^{(t)}$  and  $U^{(t)}$  denote the current estimation of  $W$  and  $U$ . When computing the stochastic gradient, instead of (2), we use<sup>1</sup>

$$\text{out}(X, A^{1t}, A^{2t}; W^{(t)}, U^{(t)}) = C\sigma(W^{(t)}XA^{1t}) + C\sigma(U^{(t)}\sigma(W^{(t)}XA^{1t})A^{2t}) \quad (5)$$

The main notations are summarized in Table 1 in Appendix.

## 3 Main Algorithm and Theoretical Results

### 3.1 Informal Key Theoretical Findings

We first summarize our major theoretical insights and takeaways before formally presenting them.

**1. The first theoretical generalization guarantee of two-hidden-layer GCNs with jumping-connection.** We demonstrate that training a single jumping-connection two-hidden-layer GCN using our Algorithm 1 returns a model that achieves the label prediction performance almost the same as the best prediction performance using a large class of target functions. We also characterize quantitatively the required number of labeled nodes, referred to as the sample complexity, to achieve the desirable prediction error. To the best of our knowledge, only Li et al. (2022a); Zhang et al. (2023b) provide explicit sample complexity bounds for node classification using graph neural networks, but for shallow GCNs with no jumping connection. Our work is the first one that provides a theoretical generalization and sample complexity analysis for the practical GCN architecture with jumping connections.

**2. Graph sparsification affects generalization through the sparse effective adjacency matrix  $A^*$ .** We show that training with edge pruning produces a model with the same prediction accuracy as a model trained on a GCN with  $A^*$  as the normalized adjacency matrix, i.e., replacing  $A$  with  $A^*$  in (1). The effective adjacency matrix is first discussed in (Li et al., 2022a), in the setup of node sampling for two-hidden-layer GCNs with no jumping connection, but  $A^*$  in (Li et al., 2022a) is dense. We show that the effective adjacency

<sup>1</sup>If different layers use different adjacency matrices, we specify both matrices in the function representation. Otherwise, we use one matrix to simplify notations.

matrix also exists for edge pruning on GCNs with jumping connection and can be sparse, indicating that the sparsified matrices can be very sparse without sacrificing generalization.

**3. Layer-wise graph sparsification due to jumping connection.** We show that in the two-hidden-layer GCN with a single jumping connection, the first hidden-layer learns a simpler base function that contributes more to the output, while the second hidden-layer learns a more complicated function that contributes less to the output. Therefore, compared with the first hidden layer, the second hidden layer is more robust to graph sparsification and can tolerate a deviation of the pruned matrix to  $A^*$  without affecting the prediction accuracy. To the best of our knowledge, this is the first theoretical characterization of how jumping connections affect the sparsification requirements in different layers, while the previous analysis in (Li et al., 2022a) assumes the same matrix sampling deviations for all layers.

### 3.2 Graph Topology Sparsification Strategy

Our theoretical sparsification strategy differs slightly from Algorithm 1 due to our adjustments aiming to facilitate and simplify the theoretical analysis. Nevertheless, our core concept is remaining more large-weight edges with a higher probability, while remaining small-weight edges with a smaller probability.

We follow the same assumption on node degrees as that in Li et al. (2022a). Specifically, the node degrees in  $\mathcal{G}$  can be divided into  $L$  ( $L \geq 1$ ) groups, with each group having  $N_l$  nodes ( $l \in [L]$ ). The degrees of all  $N_l$  nodes in group  $l$  are in the order of  $d_l$ , i.e., between  $cd_l$  and  $Cd_l$  for some constants  $c \leq C$ .  $d_l$  is order-wise smaller than  $d_{l+1}$ , i.e.,  $d_l = o(d_{l+1})$ .

Let matrix  $A_{B_{ij}} \in \mathbb{R}^{N_i \times N_j}$  denote a submatrix of  $A$  with rows in group  $i$  and columns corresponding to group  $j$ . Then all entries in  $A_{B_{ij}}$  are in the order of  $\frac{1}{\sqrt{d_i d_j}}$ . Note that a relatively smaller entry in  $A$  corresponds to an edge between relatively higher-degree nodes. Let  $p_{ij}^k$  in  $[0, 1]$  ( $k = 1, 2$ ) reflect the probability of remaining weights on smaller entries in  $A_{B_{ij}}$  (i.e., the probability of pruning weights on smaller entries in  $A_{B_{ij}}$  is  $1 - p_{ij}^k$ ) in the first and second hidden layers, respectively. Our sparsification strategy can be described as follows: at each iteration, for each submatrix  $A_{B_{ij}}$ ,

- (1) if  $i > j$ , each of the top<sup>2</sup>  $d_1 \sqrt{\frac{d_i}{d_j}}$  large-weight edges  $A_{ij}$  in  $A_{B_{ij}}$  is retained independently with probability  $1 - p_{ij}^k$ . The remaining entries in  $A_{B_{ij}}$  are retained independently with a probability of  $p_{ij}^k$ .
- (2) if  $i \leq j$ , each of the top  $d_1$  largest  $A_{ij}$  in  $A_{B_{ij}}$  are retained with probability  $1 - p_{ij}^k$ . The remaining entries in  $A_{B_{ij}}$  are retained independently with a probability of  $p_{ij}^k$ .

We allow the pruning rates to vary in different layers and will quantify how these weights affect generalization differently. When others are fixed, a small  $p_{ij}^k$  indicates retaining primarily large-weight edges, while a large  $p_{ij}^k$  indicates retaining less large-weight edges.

To see why this sparsification strategy prioritizes low-degree edges, consider the pruned edges between a fixed group  $i$  and other groups  $j$ . Assume  $p_{ij}^k$  are the same for all groups for simplicity. If  $d_j < d_{j'}$  for two groups  $j$  and  $j'$ , then  $d_1 \sqrt{\frac{d_i}{d_j}} > d_1 \sqrt{\frac{d_i}{d_{j'}}$ , indicating that more lower-degree edges (large-weight edges) connecting groups  $i$  and  $j$  are retained, compared with edges connecting groups  $i$  and  $j'$ .

To analyze the impact of this graph topology sparsification on the learning performance, we define the *sparse effective adjacency matrix*  $A^*$  where in each submatrix  $A_{B_{ij}}^*$ :

- (1) if  $i > j$ , the top  $d_1 \sqrt{\frac{d_i}{d_j}}$  largest values in  $A_{B_{ij}}$  remain the same, while other entries are set to zero.
- (2) if  $i \leq j$ , the top  $d_1$  largest values in  $A_{B_{ij}}$  remain the same, while other entries are set to zero.

One can easily check from the definition that  $\|A^*\|_1 = O(1)$ , i.e., the maximum absolute column sum of  $A^*$  is bounded by a constant. Moreover,  $A^*$  is sparse by definition.

<sup>2</sup>The values  $d_1 \sqrt{\frac{d_i}{d_j}}$  and  $d_1$  for selecting the top largest entries in  $A_{ij}$  are chosen to simplify our theoretical analysis. In fact, any values in these orders are sufficient for our theoretical analysis. Note that the main idea of retaining large-weight edges with higher probability is maintained in our sparsification strategy.

**Algorithm 1** SGD with Layer-wise Sparsification (LWS)

- 
- 1: **Input:** Graph  $\mathcal{G}$  with normalized adjacency matrix  $A$ , node features  $X$ , known labels in  $\Omega$ , step size  $\eta_w$  and  $\eta_v$ , number of iterations  $T$ , pruning rate  $p_{ij}^1$  and  $p_{ij}^2$ .
  - 2: Initialize  $W^{(0)}, V^{(0)}, C$ .  $\mathbf{W}_0 = 0, \mathbf{V}_0 = 0$
  - 3: **for**  $t = 0, 1, \dots, T - 1$  **do**
  - 4: Retain the top  $q_1$  fraction of the largest weight edges with high probability  $(1 - p_{ij}^1)$  and retain the remaining  $1 - q_1$  fraction of small weights with low probability  $(p_{ij}^1)$  to get  $A^{1t}$ .
  - 5: Retain the top  $q_2$  fraction of the largest edge weights with high probability  $(1 - p_{ij}^2)$  and retain the remaining  $1 - q_2$  fraction of small weight with low probability  $(p_{ij}^2)$  to get  $A^{2t}$ . ( $q_1 > q_2$  and  $p_{ij}^1 < p_{ij}^2$ ).
  - 6: Randomly sample  $n$  from  $\Omega$ .
  - 7: Calculate the gradient of  $L$  in (24) and update weight deviations through
 
$$\mathbf{W}_{t+1} \leftarrow \mathbf{W}_t - \eta_w \frac{\partial L(\mathbf{W}, \mathbf{V})}{\partial \mathbf{W}} \Big|_{\mathbf{W}=\mathbf{W}_t, \mathbf{V}=\mathbf{V}_t}$$

$$\mathbf{V}_{t+1} \leftarrow \mathbf{V}_t - \eta_v \frac{\partial L(\mathbf{W}, \mathbf{V})}{\partial \mathbf{V}} \Big|_{\mathbf{W}=\mathbf{W}_t, \mathbf{V}=\mathbf{V}_t}$$
  - 8: **end for**
  - 9: **Output:**  $W^{(T)} = W^{(0)} + \mathbf{W}_T, V^{(T)} = V^{(0)} + \mathbf{V}_T$ .
- 

**3.3 Concept Class and Hierarchical Learning**

In the context of GCNs, a concept class represents the set of possible target functions that map node features to labels. Defining this space is essential for understanding the function approximation capability of a learned GCN model and its ability to generalize to unseen data. Our theoretical generalization analysis establishes that the prediction error of the learned GCN model is bounded by a small constant multiple of the minimum achievable error within a well-defined concept class. This implies that the model effectively approximates the optimal function within this space. When the concept class accurately captures the true mapping from node features to labels, the minimum achievable prediction error approaches zero. Consequently, the learned GCN model also attains a low prediction error. This concept class depends on the sparsified adjacency matrix  $A^*$  rather than the original  $A$ . We show that as long as the sparsified matrices  $A^t$  remain close to  $A^*$ , even when highly sparse, graph sparsification does not degrade generalization performance.

To formally describe the concept class, consider a space of target functions  $\mathcal{H}$ , consisting of two smooth functions  $\mathcal{F}$  and  $\mathcal{G} : \mathbb{R}^{d \times N} \times \mathbb{R}^{N \times N} \rightarrow \mathbb{R}^{k \times N}$ , along with a constant  $\alpha \in \mathbb{R}^+$ :

$$\mathcal{H}_{A^*}(X) = \mathcal{F}_{A^*}(X) + \alpha \mathcal{G}_{A^*}(\mathcal{F}_{A^*}(X)), \quad (6)$$

where the  $r$ -th row ( $r \in [k]$ ) of  $\mathcal{F}_{A^*}$  and  $\mathcal{G}_{A^*}$ , denoted by  $\mathcal{F}^r$  and  $\mathcal{G}^r : \mathbb{R}^{d \times N} \times \mathbb{R}^{N \times N} \rightarrow \mathbb{R}^{1 \times N}$ , satisfy:

$$\mathcal{F}_{A^*}^r(X) = \sum_{i=1}^{p_{\mathcal{F}}} a_{\mathcal{F},r,i}^* \cdot \mathcal{F}^{r,i}(w_{r,i}^{*T} X A^*),$$

$$\mathcal{G}_{A^*}^r(X) = \sum_{i=1}^{p_{\mathcal{G}}} a_{\mathcal{G},r,i}^* \cdot \mathcal{G}^{r,i}(v_{r,i}^{*T} X A^*), \quad (7)$$

where  $p_{\mathcal{F}}, p_{\mathcal{G}}$  are the counts of basis functions used to construct the decompositions of  $\mathcal{F}^r$  and  $\mathcal{G}^r$ ;  $a_{\mathcal{F},r,i}^*, a_{\mathcal{G},r,i}^* \in [-1, 1]$  are scalar coefficients for given  $r, i$ ;  $w_{r,i}^* \in \mathbb{R}^d$  and  $v_{r,i}^* \in \mathbb{R}^k$  are vectors with norms  $\|w_{r,i}^*\| = \|v_{r,i}^*\| = \frac{1}{\sqrt{2}}$  for all  $r, i$ ;  $\mathcal{F}^{r,i}, \mathcal{G}^{r,i} : \mathbb{R} \rightarrow \mathbb{R}$  are smooth activation functions applied element-wise.

The complexities of  $\mathcal{F}$  and  $\mathcal{G}$  are represented by the tuples  $(p_{\mathcal{F}}, \mathcal{C}_s(\mathcal{F}), \mathcal{C}_\epsilon(\mathcal{F}))$  and  $(p_{\mathcal{G}}, \mathcal{C}_s(\mathcal{G}), \mathcal{C}_\epsilon(\mathcal{G}))$ , respectively.  $\mathcal{C}_\epsilon$  and  $\mathcal{C}_s$  represent model and sample complexities, respectively. The overall complexity of  $\mathcal{H}$  is quantified by the tuple  $(p_{\mathcal{F}}, p_{\mathcal{G}}, \mathcal{C}_s(\mathcal{F}), \mathcal{C}_s(\mathcal{G}), \mathcal{C}_\epsilon(\mathcal{F}), \mathcal{C}_\epsilon(\mathcal{G}))$ . The complexity of  $\mathcal{F}$  (or  $\mathcal{G}$ ) is determined by the most complex sub-target function among the  $p_{\mathcal{F}}$  (or  $p_{\mathcal{G}}$ ) smooth functions. Specifically, the complexities for

$\mathcal{F}$  and  $\mathcal{G}$  are defined as:

$$\begin{aligned} \mathcal{C}_\epsilon(\mathcal{F}) &= \max_{r,i} \{ \mathcal{C}_\epsilon(\mathcal{F}^{r,i}, \|A^*\|_1) \}, & \mathcal{C}_s(\mathcal{F}) &= \max_{r,i} \{ \mathcal{C}_s(\mathcal{F}^{r,i}, \|A^*\|_1) \}, \\ \mathcal{C}_\epsilon(\mathcal{G}) &= \max_{r,i} \{ \mathcal{C}_\epsilon(\mathcal{G}^{r,i}, \|A^*\|_1) \}, & \mathcal{C}_s(\mathcal{G}) &= \max_{r,i} \{ \mathcal{C}_s(\mathcal{G}^{r,i}, \|A^*\|_1) \}. \end{aligned} \quad (8)$$

The model and sample complexity definitions follow similarly to those in Li et al. (2022a) (Section 1.2) and Allen-Zhu & Li (2019) (Section 4). Please see Appendix B for details.

$\mathcal{F}$  and  $\mathcal{G}$  can both be viewed as one-hidden-layer GCNs with smooth activation functions and adjacency matrix  $A^*$ . The target function  $\mathcal{H}$  includes the base signal  $\mathcal{F}$ , which is less complex yet contributes significantly to the target, and  $\mathcal{G}$ , which is more complicated but contributes less. We will show that the learner networks defined in (5) can learn the concept class of target functions defined in (6). Intuitively, we will show that using a two-hidden-layer GCN with a jumping connection, the first hidden layer learns the low-complexity  $\mathcal{F}$ , and the second hidden layer learns the high-complexity  $\mathcal{G}(\mathcal{F})$  with the help of  $\mathcal{F}$  learned by the first hidden layer using the jumping connection.

We will also show that the learned GCN by our method performs almost the same as the best function in the concept class in predicting unknown labels. Let  $\mathcal{D}_{\tilde{x}_n}$  and  $\mathcal{D}_{y_n}$  denote the distribution from which the feature and label of node  $n$  are drawn, respectively. Let  $\mathcal{D}$  denote the concatenation of these distributions. Then the given feature matrix  $X$  and partial labels in  $\Omega$  can be viewed as  $|\Omega|$  identically distributed but correlated samples  $(X, y_n)$  from  $\mathcal{D}$ . The correlation results from the fact that the label of node  $i$  depends on not only the feature of node  $i$  but also neighboring features.

The  $n$ -th column of  $\mathcal{H}_{A^*}$ , denoted  $\mathcal{H}_{n,A^*} : \mathbb{R}^{d \times N} \times \mathbb{R}^{N \times N} \rightarrow \mathbb{R}^k$ , represents the target function for node  $n$ . To measure the label approximation performance of the target function, define

$$\mathbb{E}_{(X, y_n) \sim \mathcal{D}, n \in \mathcal{V}} \left[ \frac{1}{2} \|\mathcal{H}_{n,A^*}(X) - y_n\|_2^2 \right] = \text{OPT} \quad (9)$$

as the minimum prediction error achieved by the best target function in the concept class in (8). OPT decreases when the target functions are more complex, or the concept class enlarges, or if  $A^*$  characterizes the node correlations properly.

### 3.4 Modeling the prediction error of unknown labels

To simplify the analysis and representation, we re-parameterize  $U$  in (1) and (2) as  $VC$ , where  $V \in \mathbb{R}^{m \times k}$ . Then, (2) can be rewritten as follows:

$$\begin{aligned} \text{out}_n(X, A; W, V) &= \text{out}_n^1(X, A) + C\sigma(V \text{out}^1(X, A)a_n) \\ \text{where } \text{out}^1(X, A; W) &= C\sigma(WXA), \\ \text{out}_n^1(X, A; W) &= C\sigma(WXa_n) \end{aligned} \quad (10)$$

We follow the conventional setup for theoretical analysis that  $C$  is fixed at its random initialization, and only  $W$  and  $V$  are updated during training.  $C, W^{(0)}, V^{(0)}$  are randomly initialized from Gaussian distributions,  $C_{i,j} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \frac{1}{m})$ ,  $W_{i,j}^{(0)} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma_w^2)$  and  $V_{i,j}^{(0)} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma_v^2/m)$ , respectively.

The algorithm is summarized in Algorithm 1. When computing the stochastic gradient of a sampled label  $y_n$ , the loss is represented as a function of the weight deviations  $\mathbf{W}, \mathbf{V}$  from initiation  $W^{(0)}$  and  $V^{(0)}$ , i.e.,

$$L(\mathbf{W}, \mathbf{V}) = \text{Obj}_n(X, A^{1t}, A^{2t}, y_n; W^{(0)} + \mathbf{W}, V^{(0)} + \mathbf{V}). \quad (11)$$

$\mathbf{W}_t$  and  $\mathbf{V}_t$  denote the weight deviations of the estimated weights  $W^{(t)}$  and  $V^{(t)}$  in iteration  $t$  from  $W^{(0)}$  and  $V^{(0)}$ , i.e.,  $W^{(t)} = W^{(0)} + \mathbf{W}_t$ ,  $V^{(t)} = V^{(0)} + \mathbf{V}_t$ . We assume  $0 < \alpha \leq \tilde{O}\left(\frac{1}{\sigma_s(\mathcal{G})}\right)$  throughout the training. We prove that  $\|\mathbf{W}_t\|_2$  and  $\|\mathbf{V}_t\|_2$  are bounded by  $\tilde{\Theta}(\mathcal{C}_s(\mathcal{F}))$  and  $\tilde{\Theta}(\alpha\mathcal{C}_s(\mathcal{G}))$  during training, i.e.,  $\|\mathbf{W}_t\|_2 \leq \tilde{\Theta}(\mathcal{C}_s(\mathcal{F}))$ ,  $\|\mathbf{V}_t\|_2 \leq \tilde{\Theta}(\alpha\mathcal{C}_s(\mathcal{G})) < 1$  for all  $t$ .

The following lemma shows that graph sparsification in different layers contributes to the output approximation differently. In other words, to maintain the same accuracy in the output, the tolerable pruning rates in different layers are different.

**Lemma 3.1.** For any given constant  $E$ , if the first and second layer  $A^{1t}$  and  $A^{2t}$  are sparsified with  $p_{ij}^1 \leq \tilde{\Theta}(\frac{\sqrt{d_i d_j} E}{N_i N_j C_s(\mathcal{F})})$  and  $p_{ij}^2 \leq \tilde{\Theta}(\frac{\sqrt{d_i d_j} E}{N_i N_j \alpha C_s(\mathcal{F}) C_s(\mathcal{G})})$ , respectively, then with the probability over  $1 - e^{-\Omega(E\sqrt{d_i d_j}/C_s(\mathcal{F}))}$

$$\begin{aligned} \|A^{1t} - A^*\|_1 &\leq \frac{E}{\tilde{\Theta}(C_s(\mathcal{F}))}, \\ \|A^{2t} - A^*\|_1 &\leq \frac{E}{\tilde{\Theta}(\alpha C_s(\mathcal{F}) C_s(\mathcal{G}))}, \\ \|\text{out}_n(X, A^{1t}, A^{2t}; W^{(t)}, V^{(t)}) - \text{out}_n(X, A^*; W^{(t)}, V^{(t)})\|_2 &\leq E, \quad \text{for all } t. \end{aligned} \tag{12}$$

Note that  $\tilde{\Theta}(\alpha C_s(\mathcal{G})) < 1$  (see Table 3 in Appendix for the parameters), then the upper bound for  $p_{ij}^2$  is higher than that for  $p_{ij}^1$  in the assumption. That means the pruning for the first hidden layer must focus more on low-degree edges (large-weight edges), while such a requirement is relaxed in the second layer. Then, (12) indicates that a larger deviation of the sparsified matrix  $A^t$  from  $A^*$  can be tolerated in the second hidden layer compared with the first layer. Lemma 3.1 reveals that the jumping connection allows for a more flexible sparsification strategy in deeper layers.

We will show the learned model can achieve an error close to  $O(\text{OPT})$ . Our main theorem can be sketched as follows,

**Theorem 3.2.** For  $\epsilon_0 = \tilde{\Theta}(\alpha^4 C_s(\mathcal{G})^4) < 1$  and  $\epsilon = 10 \cdot \text{OPT} + \epsilon_0$ , suppose pruning probability  $p_{ij}^1 \leq \tilde{\Theta}(\frac{\sqrt{d_i d_j} \epsilon_0}{N_i N_j C_s(\mathcal{F})})$  and  $p_{ij}^2 \leq \tilde{\Theta}(\frac{\sqrt{d_i d_j} \epsilon_0}{N_i N_j \alpha C_s(\mathcal{F}) C_s(\mathcal{G})})$ , there exist  $M_0 = \text{poly}(C_s(\mathcal{F}), C_s(\mathcal{G}), \alpha^{-1})$ ,  $T_0 = \tilde{\Theta}(\frac{C_s(\mathcal{F})^2}{\|A^*\|_1 \min\{0.1, \epsilon^2\}})$  and  $N_0 = \tilde{\Theta}(\Delta^4 C_s(\mathcal{F})^2 \|A^*\|_1^4 \log N \epsilon^{-2})$  such that for every  $m \geq M_0$ ,  $T \geq T_0$  and  $|\Omega| \geq N_0$ , with high probability, the SGD algorithm satisfies

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}_{\substack{(X, y_n) \sim \mathcal{D} \\ n \in \mathcal{V}}} \|y_n - \text{out}_n(X, A^{1t}, A^{2t}; \mathbf{W}_t, \mathbf{V}_t)\|_2^2 \leq \epsilon \tag{13}$$

As a sanity check, when the concept class enlarges,  $\text{OPT}$  decreases. Theorem 3.2 shows that the required number of neurons  $M_0$  (model complexity) and labels  $N_0$  (sample complexity) both increase accordingly.  $p_{ij}^1$  and  $p_{ij}^2$  decreasing means that we should retain more large-weight edges. Thus, our theoretical bounds match the intuition that a larger model, more labels, and more high-weight edges improve the prediction accuracy.  $N_0$  is in the order of  $\log N$ , indicating that the unknown labels can be accurately predicted from partial labels. Moreover, when  $\|A^*\|_1$  increases,  $C_\epsilon$  and  $C_s$ , and  $\epsilon_0$  all increase. Theorem 3.2 indicates the model complexity  $M_0$ ,  $N_0$ , and the generalization error  $\epsilon$  all increasing, indicating worse prediction performance.

This proof of Theorem 3.2 builds upon the proof of Theorem 1 in Allen-Zhu & Li (2019), which analyzes the generalization of a three-layer ResNet for a supervised regression problem. We extend the analysis to training GCNs with graph sparsification for a semi-supervised node regression problem. The main technical challenge is to handle the dependence of labels on neighboring features and the error in adjacency matrices due to the sparsification. Compared with Li et al. (2022a) which also considers training GCN with graph sampling, we consider a different sparsification method from that in Li et al. (2022a). The resulting  $A^*$  in Li et al. (2022a) is a dense matrix as  $A$ , while  $A^*$  in our paper is a sparse matrix. Our results thus allow the sparsified matrices to be very sparse while still maintaining the generalization accuracy. Moreover, the sampling method is the same for both hidden layers in Li et al. (2022a), resulting the same deviation from  $A^*$  in both layers. Our results indicate that the jumping connection allows a more flexible sparsification method in the second layer.

### 3.5 Proof Overview

In practice, for computational efficiency, we use the sparsified adjacency matrix  $A^t$  in the learning network. Therefore, the discrepancy between the target function with  $A^*$  and the practical learning network with  $A^t$  can be viewed as two parts:

$$\begin{aligned} \|\mathcal{H}_{n, A^*}(X) - \text{out}_n(X, A^{1t}, A^{2t}; \mathbf{W}_t, \mathbf{V}_t)\|_2 &\leq \|\mathcal{H}_{n, A^*}(X) - \text{out}_n(X, A^*)\|_2 \\ &\quad + \|\text{out}_n(X, A^{1t}, A^{2t}) - \text{out}_n(X, A^*)\|_2 \end{aligned} \tag{14}$$

the first part quantifies how well the learning network, trained with  $\mathbf{W}_t$  and  $\mathbf{V}_t$  using  $A^*$ , can approximate the target function  $\mathcal{H}_{n,A^*}$ . We prove the existence of  $\mathbf{W}^*$  and  $\mathbf{V}^*$  (see Lemma C.3) with  $m \geq M_0$ ,  $T \geq T_0$  and  $\Omega \geq N_0$  such that

$$\|\mathcal{H}_{n,A^*}(X) - \text{out}_n(X, A^*; \mathbf{W}^*, \mathbf{V}^*)\|_2 \leq \epsilon_0. \quad (15)$$

The second part quantifies the difference between the learning network’s output when using the sparse adjacency matrices  $A^t$  and effective adjacency matrix  $A^*$ . Specifically, it is represented by the term:

$$\begin{aligned} \|\text{out}_n(X, A^{1t}, A^{2t}) - \text{out}_n(X, A^*)\|_2 &\leq \|C\sigma(WXa_n^{1t}) - C\sigma(WXa_n^*)\|_2 + \\ &\|C\sigma(V\text{out}_n^1(XA^{1t})a_n^{2t}) - C\sigma(V\text{out}_n^1(XA^*)a_n^*)\|_2 \end{aligned} \quad (16)$$

For the inequality  $\|C\sigma(WXa_n^{1t}) - C\sigma(WXa_n^*)\|_2 \leq \epsilon_0$  to hold, it is required that  $\|A^{1t} - A^*\|_1 \leq \frac{\epsilon_0}{\Theta(C_s(\mathcal{F}))}$  and similarly,  $\|A^{2t} - A^*\|_1 \leq \frac{\epsilon_0}{\Theta(\alpha C_s(\mathcal{F})C_s(\mathcal{G}))}$  (see Appendix C.4). We establish that with appropriate pruning probabilities  $p_{ij}^1$  and  $p_{ij}^2$ , the norms  $\|A^{1t} - A^*\|_1$  and  $\|A^{2t} - A^*\|_1$  can be sufficiently small (see Appendix C.7).

Finally, consider the definition of OPT, we can prove  $\|y_n - \text{out}_n(X, A^{1t}, A^{2t}; \mathbf{W}_t, \mathbf{V}_t)\|_2 \leq \epsilon$ .

## 4 Empirical Experiment

### 4.1 Experiment on synthetic data

We generate a graph  $\mathcal{G}$  with  $N = 2000$  nodes. Given  $A$ , the sparse  $A^*$  is obtained following the procedure in Section 3.2. The node labels are generated by the target function

$$\begin{aligned} \mathcal{F}_{A^*}(X) &= CW^*XA^*, \\ \mathcal{G}_{A^*}(\mathcal{F}_{A^*}(X)) &= C(\sin(V^*\mathcal{F}_{A^*}(X)A^*) \odot \tanh(V^*\mathcal{F}_{A^*}(X)A^*)), \\ \mathcal{H}_{A^*}(X) &= \mathcal{F}_{A^*}(X) + \alpha\mathcal{G}_{A^*}(\mathcal{F}_{A^*}(X)). \end{aligned} \quad (17)$$

where  $X \in \mathbb{R}^{d \times N}$ ,  $W^* \in \mathbb{R}^{r \times d}$ ,  $V^* \in \mathbb{R}^{r \times k}$ , and  $C \in \mathbb{R}^{k \times r}$  are randomly generated with each entry i.i.d. from  $\mathcal{N}(0, 1)$ .  $d = 100$ ,  $k = 5$ ,  $r = 30$ ,  $\alpha = 0.5$ . A two-hidden-layer GCN with a single jumping connection as defined in (2) with  $m$  neurons in each hidden layer is trained on a randomly selected set  $\Omega$  of labeled nodes. The rest  $N - |\Omega|$  labels are used for testing. The test error is measured by the  $\ell_2$  regression loss in (3).

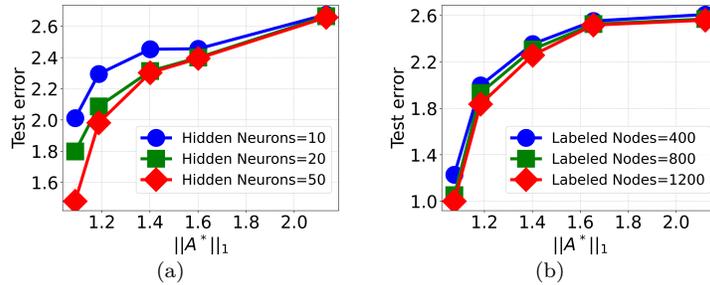


Figure 1: Experiment on two-degree group synthetic data: (a) Test error with  $|\Omega| = 1200$ . (b) Test error with  $m = 50$ .

**Model and sample complexities with  $\|A^*\|_1$ :** In Figures 1,  $\mathcal{G}$  has two-degree groups. Group 1 has  $N_1 = 200$  nodes, and the degree of each node follows a Gaussian distribution  $\mathcal{N}(d_1, \sigma^2)$ . Group 2 has  $N_2 = 1800$  nodes, and the degree of each node follows a Gaussian distribution  $\mathcal{N}(d_2, \sigma^2)$ . The degrees are truncated to fall within the range of 0 to 500. We vary  $A^*$  by changing  $d_2$  and the corresponding  $A$ . We fix  $d_1 = 200$  and  $\sigma = 20$ . We directly train with  $A^*$  to study the impact of  $A^*$  on model and sample complexities. In Figures 1 (a),  $|\Omega| = 1200$  and the number of neurons per layer  $m$  varies. To achieve the same test accuracy, when  $\|A^*\|_1$  increases, the number of neurons also increases, verifying our model complexity  $M_0$  in Theorem 3.2. In Figures 1 (b),  $m = 50$  and  $|\Omega|$  varies. To achieve the same test accuracy, when  $\|A^*\|_1$  increases, the required number of labels also increases, verifying our sample complexity  $N_0$  in Theorem 3.2.

In Figures 2,  $\mathcal{G}$  has one-degree groups and the degree of each node follows a Gaussian distribution  $\mathcal{N}(d, \sigma^2)$ .  $d = 200$ . The degrees are truncated to fall within the range of 0 to 500. We vary  $A^*$  by changing  $\sigma$  and the corresponding  $A$ . We directly train with  $A^*$  to study the impact of  $A^*$  on model and sample complexities. In Figures 2 (a),  $|\Omega| = 1200$  and the number of neurons per layer  $m$  varies. Fig. 2 (a) shows the testing error decreases as  $m$  increases. When  $m$  is the same, the testing error increases as  $\|A^*\|_1$  increases. This verifies our model complexity in Theorem 3.2. In Figures 2 (b),  $m = 50$  and  $|\Omega|$  varies. Fig. 2 (b) shows the testing error decreases as  $\Omega$  increases. When  $\Omega$  is the same, the testing error increases as  $\|A^*\|_1$  increases. This verifies our model complexity in Theorem 3.2.

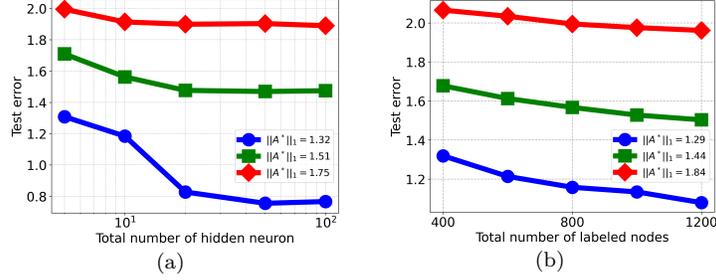


Figure 2: Experiment on one-degree group synthetic data: (a) Test error with  $|\Omega| = 1200$ . (b) Test error with  $m = 50$ .

**Layer-wise Sparsification impact on generalization:** We fix  $\Omega = 1200$ ,  $m = 50$ ,  $\|A^*\|_1 = 1.27$ . We sample adjacency matrices during training. Figure 3 shows the relationship between the test error and the average deviation of sparsified matrices ( $A^{1t}$  and  $A^{2t}$  in the first and second hidden layers) from  $A^*$ . We can see that pruning in the second hidden layer (yellow dashed line) contributes to generalization degradation much milder than pruning in the first hidden layer (red solid arrow). This verifies our Lemma 3.1 that the sparsification requirements are more restrictive in the first layer than the second layer to maintain the same generalization accuracy.

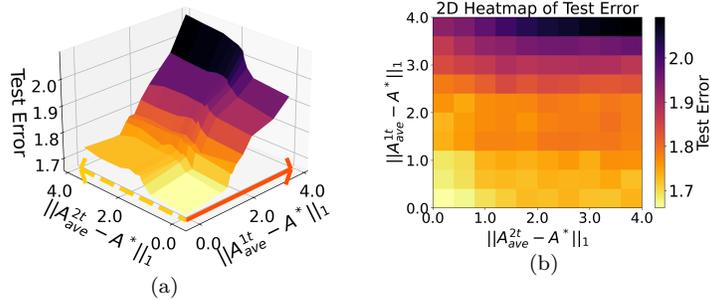


Figure 3: Experiment on synthetic data of layer-wise sparsification. (a) The second hidden layer tolerates more test error than the first hidden layer. (b) 2D heatmap of test error.

## 4.2 Experiment on real data

**Retaining large weights in the graph can perform as well as trained sparse graph methods.** We applied Algorithm 1 to a one-hidden-layer shallow GCN on small-scale datasets (Cora, Citeseer) for node multi-class classification tasks, comparing it with state-of-the-art (SOTA) sparsification methods such as CGP (Liu et al., 2023) and UGS (Chen et al., 2021). For these small datasets, multi-layer GCNs are not necessary, so we employ the one-hidden-layer GCN (Kipf & Welling, 2017) with 512 hidden neurons and preclude the use of different pruning rates for shallow and deep layers. We retain the top  $q$  fraction of the largest edge weights with a 99% probability and retain the remaining  $1 - q$  fraction of small weight with a 1% probability to get  $A^t$ , so the sparsify of our method (LWS) is  $0.98q + 0.01$  and we vary  $q$  from 0.01 to 1.0. In Figure 4, we only demonstrate that retaining large-weight edges from  $A$  using our method can achieve performance comparable to that of trained sparse graphs produced by SOTA methods.

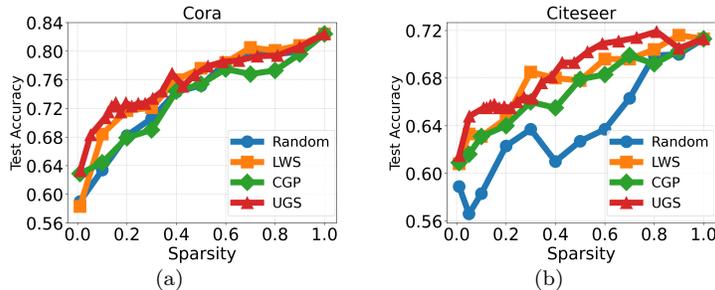


Figure 4: Experiment on sparsifying shallow GCN models.

We also evaluate multi-layer GCNs with jumping connections on the large-scale Open Graph Benchmark (OGB) datasets for node multi-class classification tasks. A summary of Ogbn datasets’ statistics is presented in Table 2 in Appendix.

The task of Ogbn-Arxiv is to classify the 40 subject areas of arXiv CS papers. We use 60% of the data for training, 20% for testing, and 20% for verification. We deploy an 8-layer Jumping Knowledge Network (Xu et al., 2018b) GCN with concatenation layer aggregation as a learner network. We treat the first four layers as shallow layers and the last four layers as deep layers. Shallow and deep layers are sparsified differently. The generalization is evaluated by the fraction of erroneous predictions of unknown labels.

**Pruning in deep layers is more flexible with less generalization degradation.** In this experiment, we employ a simplified version of the graph sparsification method discussed in Section 3.2. For the shallow layers, at each iteration  $t$ , we obtain a sparsified adjacency matrix  $A^{1t}$  as follows: we retain the top  $q_1$  fraction of the largest weight edges  $A_{ij}$  from the adjacency matrix  $A$  with a 99% probability, and retain the remaining entries with a 1% probability. For the deep layers, the sparsified adjacency matrix  $A^{2t}$  is generated similarly, but we use the top  $q_2$  fraction of largest  $A_{ij}$ , again retaining with probabilities of 99% and 1% for the top and remaining entries, respectively.

Figure 5 shows test error when  $q_1$  and  $q_2$  vary. One can see that the test error decreases more drastically when only increasing  $q_1$  (yellow dashed arrow) compared with only increasing  $q_2$  (red solid arrow), indicating that graph pruning in shallow layers has a more significant impact than graph pruning in deeper layers. When both  $q_1$  and  $q_2$  are greater than 0.6, the test error is always small (less than 0.29) for a wide range of  $q_1, q_2$ . That may suggest the existence of multiple sparse  $A^*$  such that sparsified matrices  $A^t$  with different  $q_1, q_2$  pairs approximate different  $A^*$ , and all  $A^*$  can accurately represent the data correlations in the mapping function from the features to the labels.

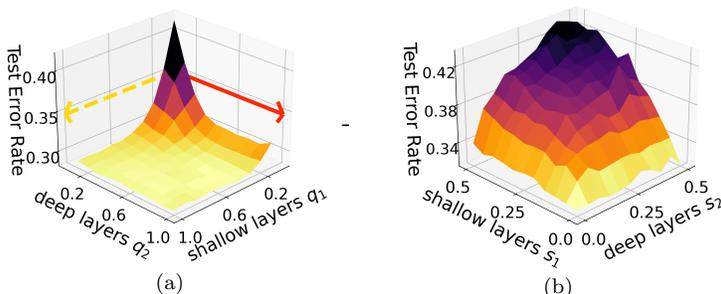


Figure 5: Learning deep GCNs on Ogbn-Arxiv: (a) Deeper layers tolerate higher sampling rates than shallow layers while maintaining accuracy. (b) 2D heatmap of test error rate.

**Large-weight edges are more influential on generalization than small-weight edges.** Note that if nodes  $i$  and  $j$  have higher degrees, then  $A_{ij}$  has a smaller value. We sparsify one matrix for the shallow

layers by keeping the values of  $A_{ij}$  that are in the range of top  $s_1$  to  $s_1 + 0.5$  fraction and setting all other values to zero. Similarly, we sparsify one matrix for the deep layers by keeping the values in the range of top  $s_2$  to  $s_2 + 0.5$  fraction and setting all other values to zero. These two sparsified matrices are used during training. When  $s_1$  and  $s_2$  increase, the resulting matrices have the same number of nonzero entries, and the sparsified entries focus more on high-degree edges. Figure 6 shows the test error indeed increases as  $s_1, s_2$  increases. This justifies the sparsification strategy to retain more large-weight edges.

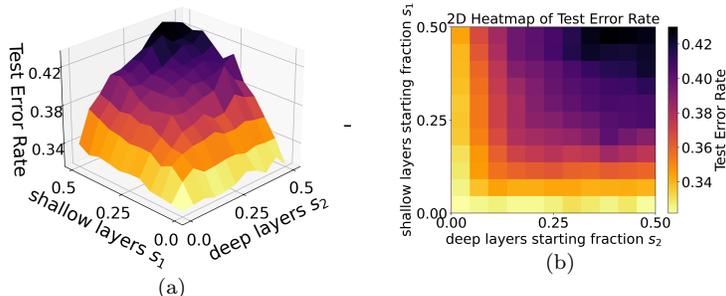


Figure 6: Learning deep GCNs on Ogbn-Arxiv: (a) Retaining more large-weight edges (small  $s_1, s_2$ ) outperforms retaining more small-weight edges (large  $s_1, s_2$ ). (b) 2D heatmap of test error rate.

For the Ogbn-Products dataset, we deploy a 4-layer Jumping Knowledge Network (Xu et al., 2018b) GCN with concatenation layer aggregation. Each hidden layer consists of 128 neurons. We define the first two layers as shallow layers and the last two layers as deep layers with sampling rate  $p_2$ . The task is to classify the category of a product in a multi-class, where the 47 top-level categories are used for target labels. We use 60% of the data for training, 20% for testing, and 20% for verification.

We run the similar experiment as Figure 5 for Ogbn-Products dataset. We fix  $q_1 = 0.1$  and vary  $q_2$  from 0.1 to 1.0 at increments of 0.1, then fix  $q_2 = 0.1$  and vary  $q_1$  from 0.1 to 1.0 at increments of 0.1. Figure 7 shows that with the increasing sampling rate in shallow layers, the test accuracy is higher than the test accuracy with the increasing sampling rate in deep layers. It suggests that the generalization is more sensitive to the sampling in the shallow layers rather than deep layers, consistent with observations in other datasets.

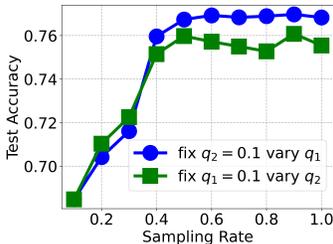


Figure 7: Layer-wise Sampling Rate Effect on Ogbn-Products

## 5 Conclusion

This paper provides a theoretical generalization analysis of training GCNs with skip connections using graph sampling. To the best of our knowledge, this paper provides the first analysis of how skip connection affects the generalization performance. We show that for a two-hidden-layer GCN with a skip connection, the first hidden layer learns a simpler function that contributes significantly to the output, making the choice of sampling more crucial in the first hidden layer. In contrast, the second layer learns a composite function that contributes less to the output, allowing for a more flexible sampling approach while preserving generalization. This insight is verified on deep GCNs on benchmark datasets. Future works include generalization analysis of deep GCNs and designing layer-specific sampling strategies that optimize generalization within the constraints of available computational resources.

## References

- Zeyuan Allen-Zhu and Yuanzhi Li. What can resnet learn efficiently, going beyond kernels? *Advances in Neural Information Processing Systems*, 32, 2019.
- Zeyuan Allen-Zhu and Yuanzhi Li. Feature purification: How adversarial training performs robust deep learning. In *2021 IEEE 62nd Annual Symposium on Foundations of Computer Science (FOCS)*, pp. 977–988. IEEE, 2022.
- Zeyuan Allen-Zhu and Yuanzhi Li. Towards understanding ensemble, knowledge distillation and self-distillation in deep learning. In *The Eleventh International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=Uuf2q9TfXGA>.
- Jianfei Chen, Tianyi Ma, and Cao Xiao. Fastgcn: Fast learning with graph convolutional networks via importance sampling. In *International Conference on Learning Representations (ICLR)*, 2018.
- Ting Chen, Yidan Sui, Xiaohan Chen, Anima Zhang, and Zhihua Wang. A unified lottery ticket hypothesis for graph neural networks. In *International Conference on Machine Learning*, pp. 1695–1706. PMLR, 2021.
- Zhengdao Chen, Soledad Villar, Lei Chen, and Joan Bruna. On the equivalence between graph isomorphism testing and function approximation with gnns. *Advances in neural information processing systems*, 32, 2019.
- Wei-Lin Chiang, Xuanqing Liu, Si Si, Yang Li, Samy Bengio, and Cho-Jui Hsieh. Cluster-gcn: An efficient algorithm for training deep and large graph convolutional networks. In *Proceedings of the 25th ACM SIGKDD international conference on knowledge discovery & data mining*, pp. 257–266, 2019.
- Weilin Cong, Morteza Ramezani, and Mehrdad Mahdavi. On provable benefits of depth in training graph convolutional networks. *Advances in Neural Information Processing Systems*, 34:9936–9949, 2021.
- Simon S. Du, Keyulu Hou, Ruslan R. Salakhutdinov, Barnabás Póczos, Ruosong Wang, and Keyulu Xu. Graph neural tangent kernel: Fusing graph neural networks with graph kernels. In *Advances in Neural Information Processing Systems*, pp. 5724–5734, 2019a.
- Simon S. Du, Xiyu Zhai, Barnabás Póczos, and Aarti Singh. Gradient descent provably optimizes over-parameterized neural networks. In *International Conference on Learning Representations (ICLR)*, 2019b.
- Keyu Duan, Zirui Liu, Peihao Wang, Wenqing Zheng, Kaixiong Zhou, Tianlong Chen, Xia Hu, and Zhangyang Wang. A comprehensive study on large-scale graph training: Benchmarking and rethinking. *Advances in Neural Information Processing Systems*, 35:5376–5389, 2022.
- Pascal Esser, Leena Chennuru Vankadara, and Debarghya Ghoshdastidar. Learning theory can (sometimes) explain generalisation in graph neural networks. *Advances in Neural Information Processing Systems*, 34: 27043–27056, 2021.
- Jiarui Feng, Yixin Chen, Fuhai Li, Anindya Sarkar, and Muhan Zhang. How powerful are k-hop message passing graph neural networks. *Advances in Neural Information Processing Systems*, 35:4776–4790, 2022.
- William L. Hamilton, Zhitao Ying, and Jure Leskovec. Inductive representation learning on large graphs. In *Advances in Neural Information Processing Systems*, pp. 1024–1034, 2017.
- Wenbing Huang, Tong Zhang, Yu Rong, and Junzhou Huang. Adaptive sampling towards fast graph representation learning. In *Advances in Neural Information Processing Systems*, pp. 4558–4567, 2018.
- Vassilis N Ioannidis, Siheng Chen, and Georgios B Giannakis. Pruned graph scattering transforms. In *International Conference on Learning Representations*, 2020.
- Arthur Jacot, Franck Gabriel, and Clement Hongler. Neural tangent kernel: Convergence and generalization in neural networks. In *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc., 2018.

- Thomas N. Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. In *International Conference on Learning (ICLR)*, 2017.
- Guohao Li, Matthias Muller, Ali Thabet, and Bernard Ghanem. Deepgens: Can gens go as deep as cnns? In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 9267–9276, 2019.
- Hongkang Li, Meng Wang, Sijia Liu, Pin-Yu Chen, and Jinjun Xiong. Generalization guarantee of training graph convolutional networks with graph topology sampling. In *International Conference on Machine Learning (ICML)*, pp. 13014–13051. PMLR, 2022a.
- Hongkang Li, Shuai Zhang, and Meng Wang. Learning and generalization of one-hidden-layer neural networks, going beyond standard gaussian data. In *2022 56th Annual Conference on Information Sciences and Systems (CISS)*, pp. 37–42. IEEE, 2022b.
- Jundong Li, Ting Zhang, Hanghang Tian, Shengnan Jin, Mingyi Fardad, and Reza Zafarani. Sgcn: A graph sparsifier based on graph convolutional networks. In *Pacific-Asia Conference on Knowledge Discovery and Data Mining*, pp. 275–287. Springer, 2020.
- Qimai Li, Zhichao Han, and Xiao-Ming Wu. Deeper insights into graph convolutional networks for semi-supervised learning. In *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.
- Yuanzhi Li and Yingyu Liang. Learning overparameterized neural networks via stochastic gradient descent on structured data. *Advances in neural information processing systems*, 31, 2018.
- Renjie Liao, Raquel Urtasun, and Richard Zemel. A pac-bayesian approach to generalization bounds for graph neural networks. In *International Conference on Learning Representations*, 2020.
- Chuang Liu, Xueqi Ma, Yibing Zhan, Liang Ding, Dapeng Tao, Bo Du, Wenbin Hu, and Danilo P. Mandic. Comprehensive graph gradual pruning for sparse training in graph neural networks. *IEEE Transactions on Neural Networks and Learning Systems*, 2023.
- Christopher Morris, Martin Ritzert, Matthias Fey, William L Hamilton, Jan Eric Lenssen, Gaurav Rattan, and Martin Grohe. Weisfeiler and leman go neural: Higher-order graph neural networks. In *Proceedings of the AAAI conference on artificial intelligence*, 2019.
- Will Norcliffe-Brown, Stathis Vafeias, and Sarah Parisot. Learning conditioned graph structures for interpretable visual question answering. *Advances in neural information processing systems*, 31, 2018.
- Kenta Oono and Taiji Suzuki. Graph neural networks exponentially lose expressive power for node classification. *arXiv preprint arXiv:1905.10947*, 2019.
- Kenta Oono and Taiji Suzuki. Optimization and generalization analysis of transduction through gradient boosting and application to multi-scale graph neural networks. *Advances in Neural Information Processing Systems*, 33:18917–18930, 2020.
- Lianke Qin, Zhao Song, and Baocheng Sun. Is solving graph neural tangent kernel equivalent to training graph neural network? *arXiv preprint arXiv:2309.07452*, 2023.
- Morteza Ramezani, Weilin Cong, Mehrdad Mahdavi, Anand Sivasubramaniam, and Mahmut Kandemir. Gcn meets gpu: Decoupling “when to sample” from “how to sample”. *Advances in Neural Information Processing Systems*, 33:18482–18492, 2020.
- Rakshith S Srinivasa, Cao Xiao, Lucas Glass, Justin Romberg, and Jimeng Sun. Fast graph attention networks using effective resistance based graph sparsification. *arXiv preprint arXiv:2006.08796*, 2020.
- Yiding Sui, Xiaojie Wang, Ting Chen, Xiaoqiang He, and Tat-Seng Chua. Inductive lottery ticket learning for graph neural networks. In *International Conference on Learning Representations (ICLR)*, pp. 1–18, 2022.

- Huayi Tang and Yong Liu. Towards understanding the generalization of graph neural networks. In *Fortieth International Conference on Machine Learning (ICML)*, 2023a.
- Huayi Tang and Yong Liu. Towards understanding generalization of graph neural networks. In *International Conference on Machine Learning*, pp. 33674–33719. PMLR, 2023b.
- Xiyuan Wang and Muhan Zhang. How powerful are spectral graph neural networks. In *International Conference on Machine Learning*, pp. 23341–23362. PMLR, 2022.
- Qitian Wu, Hengrui Zhang, Xiaofeng Gao, Peng He, Paul Weng, Han Gao, and Guihai Chen. Dual graph attention networks for deep latent representation of multifaceted social effects in recommender systems. In *The world wide web conference*, pp. 2091–2102, 2019.
- Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How powerful are graph neural networks? In *International Conference on Learning Representations*, 2018a.
- Keyulu Xu, Chengtao Li, Yonglong Tian, Tomohiro Sonobe, Ken-ichi Kawarabayashi, and Stefanie Jegelka. Representation learning on graphs with jumping knowledge networks. In *International conference on machine learning*, pp. 5453–5462. PMLR, 2018b.
- Keyulu Xu, Wei Hu, Jure Leskovec, and Stefanie Jegelka. How powerful are graph neural networks? In *International Conference on Learning Representations (ICLR)*, 2019.
- Keyulu Xu, Mengshi Zhang, Stefanie Jegelka, and Kenji Kawaguchi. Optimization of graph neural networks: Implicit acceleration by skip connections and more depth. In *International Conference on Machine Learning (ICML)*. PMLR, 2021.
- Chenxiao Yang, Qitian Wu, David Wipf, Ruoyu Sun, and Junchi Yan. How graph neural networks learn: Lessons from training dynamics. In *Forty-first International Conference on Machine Learning*, 2023.
- Rex Ying, Ruining He, Kaifeng Chen, Pong Eksombatchai, William L Hamilton, and Jure Leskovec. 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.
- Haoxiang You, Changxiao Li, Pengcheng Xu, Yonggan Fu, Yang Wang, Xiaohan Chen, Richard G. Baraniuk, Zhihua Wang, and Yingyan Lin. Drawing early-bird tickets: Toward more efficient training of deep networks. In *International Conference on Learning Representations (ICLR)*, 2020.
- Hongyang Zeng, Hongyuan Zhou, Abhay Srivastava, Rajgopal Kannan, and Viktor Prasanna. Graphsaint: Graph sampling based inductive learning method. In *International Conference on Learning Representations (ICLR)*, 2020.
- Bohang Zhang, Shengjie Luo, Liwei Wang, and Di He. Rethinking the expressive power of gnns via graph biconnectivity. In *The Eleventh International Conference on Learning Representations*, 2023a.
- Shuai Zhang, Meng Wang, Pin-Yu Chen, Sijia Liu, Songtao Lu, and Miao Liu. Joint edge-model sparse learning is provably efficient for graph neural networks. In *The Eleventh International Conference on Learning Representations*, 2023b. URL [https://openreview.net/forum?id=4UldFtZ\\_CVF](https://openreview.net/forum?id=4UldFtZ_CVF).
- Sijia Zhang, Mengqi Wang, Shichao Liu, Pin-Yu Chen, and Jinjun Xiong. Fast learning of graph neural networks with guaranteed generalizability: One-hidden-layer case. In *International Conference on Machine Learning*, pp. 11268–11277, 2020.
- Wentao Zhang, Ziqi Yin, Zeang Sheng, Yang Li, Wen Ouyang, Xiaosen Li, Yangyu Tao, Zhi Yang, and Bin Cui. Graph attention multi-layer perceptron. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pp. 4560–4570, 2022.
- Yuhao Zhang, Peng Qi, and Christopher D Manning. Graph convolution over pruned dependency trees improves relation extraction. *arXiv preprint arXiv:1809.10185*, 2018.

- Kai Zhong, Zhao Song, Prateek Jain, Peter L Bartlett, and Inderjit S Dhillon. Recovery guarantees for one-hidden-layer neural networks. In *International conference on machine learning*, pp. 4140–4149. PMLR, 2017.
- D. Zou, Z. Hu, Y. Wang, S. Jiang, Y. Sun, and Q. Gu. Layer-dependent importance sampling for training deep and large graph convolutional networks. In *Advances in Neural Information Processing Systems*, pp. 11249–11259, 2019.

## A Preliminaries

We first restate some important notations used in the Appendix, which are summarized in Table 1.

Table 1: Summary of Notations

Notations	Annotation
$\mathcal{G} = \{\mathcal{V}, \mathcal{E}\}$	$\mathcal{G}$ is an undirected graph consisting of a set of nodes $\mathcal{V}$ and a set of edges $\mathcal{E}$ .
$N$	The total number of nodes in a graph.
$A = D^{-\frac{1}{2}} \tilde{A} D^{-\frac{1}{2}}$	$A \in \mathbb{R}^{N \times N}$ is the normalized adjacency matrix computed by the degree matrix $D$ and the initial adjacency matrix $\tilde{A}$ .
$A^*$	The effective adjacency matrix.
$A^{1t}, A^{2t}$	The sparsified adjacency matrices $A$ in the first and second hidden layers at the $t$ -th iteration, respectively.
$M_0$	The required number of neurons (model complexity).
$T_0$	The required number of iterations for convergence in the SGD algorithm.
$N_0$	The required number of labeled samples (sample complexity).
$W \in \mathbb{R}^{m \times d}$	The weight matrix for the first hidden layer.
$U \in \mathbb{R}^{m \times m}$	The weight matrix for the second hidden layer.
$C \in \mathbb{R}^{k \times m}$	The weight matrix for the output layer.
$V \in \mathbb{R}^{m \times k}$	The re-parameterized weight matrix used in place of $U$ in the second hidden layer.
$a_n \in \mathbb{R}^N$	The $n$ -th column of the adjacency matrix $A$ , representing the connectivity of node $n$ .

A summary of Ognb datasets' statistics in Table 2:

Table 2: Transposed Ognb datasets statistics.

Dataset	Ognb-Arxiv	Dataset	Ognb-Products
<b>Nodes</b>	169,343	<b>Nodes</b>	2,449,029
<b>Edges</b>	1,166,243	<b>Edges</b>	61,859,140
<b>Features</b>	128	<b>Features</b>	100
<b>Classes</b>	40	<b>Classes</b>	47
<b>Metric</b>	Accuracy	<b>Metric</b>	Accuracy

**Lemma A.1.** *If  $M \in \mathbb{R}^{n \times m}$  is a random matrix where  $M_{i,j}$  are i.i.d. from  $\mathcal{N}(0, 1)$ . Then,*

- For any  $t \geq 1$ , with probability  $1 - e^{-t^2}$ , it satisfies

$$\|M\|_2 \leq O(\sqrt{n} + \sqrt{m}) + t.$$

- If  $1 \leq s \leq O\left(\frac{m}{\log^2 m}\right)$ , then with probability  $1 - e^{-(n+s \log^2 m)}$  it satisfies

$$\|Mv\|_2 \leq O\left(\sqrt{n} + \sqrt{s \log m}\right) \|v\|_2$$

for all  $s$ -sparse vectors  $v \in \mathbb{R}^m$ .

Proof: The statement can be found in Proposition B.2. from Allen-Zhu & Li (2019)

**Lemma A.2.** *Suppose  $\delta \in [0, 1]$  and  $g^{(0)} \in \mathbb{R}^m$  is a random vector  $g^{(0)} \sim \mathcal{N}(0, I_m)$ . With probability at least  $1 - e^{-\Omega(m\delta^{2/3})}$ , for all vectors  $g' \in \mathbb{R}^m$  with  $\|g'\|_2 \leq \delta$ , letting  $D' \in \mathbb{R}^{m \times m}$  be the diagonal matrix where  $(D')_{k,k} = \mathbf{1}_{(g^{(0)}+g')_k} - \mathbf{1}_{(g^{(0)})_k}$  for each  $k \in [m]$ , we have*

$$\|D'\|_0 \leq O(m^{2/3}) \quad \text{and} \quad \|D'g^{(0)}\|_2 \leq \|g'\|_2.$$

Proof: The statement can be found in Proposition B.4. from Allen-Zhu & Li (2019)

**Lemma A.3.** *Given a sampling set  $X = \{x_n\}_{n=1}^N$  that contains  $N$  partly dependent random variables, for each  $n \in [N]$ , suppose  $x_n$  is dependent with at most  $d_X$  random variables in  $X$  (including  $x_n$  itself), and the moment generating function of  $x_n$  satisfies  $\mathbb{E}[e^{sx_n}] \leq e^{Cs^2}$  for some constant  $C$  that may depend on  $x_n$ . Then, the moment generation function of  $\sum_{n=1}^N x_n$  is bounded as*

$$\mathbb{E}[e^{s \sum_{n=1}^N x_n}] \leq e^{Cd_X N s^2}.$$

Proof: The statement can be found in Lemma 7 from Zhang et al. (2020)

**Lemma A.4.**  $\|Xa_n\| \leq \|A\|_1$ .

Proof:

$$\begin{aligned} \|Xa_n\| &= \left\| \sum_{k=1}^N x_k a_{k,n} \right\| \\ &= \left\| \sum_{k=1}^N \frac{a_{k,n}}{\sum_{k=1}^N a_{k,n}} x_k \right\| \cdot \sum_{k=1}^N a_{k,n} \\ &\leq \sum_{k=1}^N \frac{a_{k,n}}{\sum_{k=1}^N a_{k,n}} \|x_k\| \cdot \|A\|_1 \\ &= \|A\|_1 \end{aligned} \tag{18}$$

**Lemma A.5.** *If  $\mathcal{F} : \mathbb{R}^d \rightarrow \mathbb{R}^k$  has general complexity  $(p, \mathfrak{C}_{\mathfrak{s}}(\mathcal{F}), \mathfrak{C}_{\varepsilon}(\mathcal{F}))$ , then for every  $x, y \in \mathbb{R}^d$ , it satisfies  $\|\mathcal{F}(x)\|_2 \leq \sqrt{k}p\mathfrak{C}_{\mathfrak{s}}(\mathcal{F}) \cdot \|x\|_2$  and  $\|\mathcal{F}(x) - \mathcal{F}(y)\|_2 \leq \sqrt{k}p\mathfrak{C}_{\mathfrak{s}}(\mathcal{F}) \cdot \|x - y\|_2$ .*

Proof: The boundedness of  $\|\mathcal{F}(x)\|_2$  is trivial so we only focus on  $\|\mathcal{F}(x) - \mathcal{F}(y)\|_2$ . For each component  $g(x) = \mathcal{F}_{r,i} \left( \frac{\langle w_{1,i}^*(x,1) \rangle}{\|(x,1)\|_2} \right) \cdot \langle w_{2,i}^*(x,1) \rangle$ , denoting by  $w_1^*$  as the first  $d$  coordinate of  $w_{1,i}^*$ , and by  $w_{2,i}^*$  as the first  $d$  coordinates of  $w_{2,i}^*$ , we have

$$\begin{aligned} g'(x) &= \mathcal{F}_{r,i} \left( \frac{\langle w_{1,i}^*(x,1) \rangle}{\|(x,1)\|_2} \right) \cdot w_{2,i}^* \\ &\quad + \langle w_{2,i}^*(x,1) \rangle \cdot \mathcal{F}'_{r,i} \left( \frac{\langle w_{1,i}^*(x,1) \rangle}{\|(x,1)\|_2} \right) \cdot \frac{w_1^* \cdot \|(x,1)\|_2 - \langle w_{1,i}^*(x,1) \rangle \cdot (x,1) / \|(x,1)\|_2^2}{\|(x,1)\|_2^2} \end{aligned}$$

This implies

$$\|g'(x)\|_2 \leq \left| \mathcal{F}_{r,i} \left( \frac{\langle w_{1,i}^*(x,1) \rangle}{\|(x,1)\|_2} \right) \right| + 2 \left| \mathcal{F}'_{r,i} \left( \frac{\langle w_{1,i}^*(x,1) \rangle}{\|(x,1)\|_2} \right) \right| \leq 3\mathfrak{C}_{\mathfrak{s}}(\mathcal{F}_{r,i})$$

As a result,  $|\mathcal{F}_r(x) - \mathcal{F}_r(y)| \leq 3p\mathfrak{C}_{\mathfrak{s}}(\mathcal{F}_{r,i})$ .

**Lemma A.6.** *For every smooth function  $\phi$ , every  $\epsilon \in (0, \frac{1}{\mathcal{C}(\phi,a)\sqrt{a^2+1}})$ , there exists a function  $h : \mathbb{R}^2 \rightarrow [-\mathcal{C}_\epsilon(\phi,a)\sqrt{a^2+1}, \mathcal{C}_\epsilon(\phi,a)\sqrt{a^2+1}]$  that is also  $\mathcal{C}_\epsilon(\phi,a)\sqrt{a^2+1}$ -Lipschitz continuous on its first coordinate with the following two (equivalent) properties:*

(a) For every  $x_1 \in [-a, a]$  where  $a > 0$ :

$$\left| \mathbb{E} \left[ \mathbf{1}_{\alpha_1 x_1 + \beta_1 \sqrt{a^2 - x_1^2} + b_0 \geq 0} h(\alpha_1, b_0) \right] - \phi(x_1) \right| \leq \epsilon$$

where  $\alpha_1, \beta_1, b_0 \sim \mathcal{N}(0, 1)$  are independent random variables.

(b) For every  $\mathbf{w}^*, \mathbf{x} \in \mathbb{R}^d$  with  $\|\mathbf{w}^*\|_2 = 1$  and  $\|\mathbf{x}\| \leq a$ :

$$\left| \mathbb{E} \left[ \mathbf{1}_{\mathbf{w}^T \mathbf{x} + b_0 \geq 0} h(\mathbf{w}^T \mathbf{w}^*, b_0) \right] - \phi(\mathbf{w}^{*\top} \mathbf{x}) \right| \leq \epsilon$$

where  $\mathbf{w} \sim \mathcal{N}(0, \mathbf{I})$  is an  $d$ -dimensional Gaussian,  $b_0 \sim \mathcal{N}(0, 1)$ .

Furthermore, we have  $\mathbb{E}_{\alpha_1, b_0 \sim \mathcal{N}(0, 1)} [h(\alpha_1, b_0)^2] \leq (C_s(\phi, a))^2 (a^2 + 1)$ .

(c) For every  $\mathbf{w}^*, \mathbf{x} \in \mathbb{R}^d$  with  $\|\mathbf{w}^*\|_2 = 1$ , let  $\tilde{\mathbf{w}} = (\mathbf{w}, b_0) \in \mathbb{R}^{d+1}$ ,  $\tilde{\mathbf{x}} = (\mathbf{x}, 1) \in \mathbb{R}^{d+1}$  with  $\|\tilde{\mathbf{x}}\| \leq \sqrt{a^2 + 1}$ , then we have

$$\left| \mathbb{E} \left[ \mathbf{1}_{\tilde{\mathbf{w}}^\top \tilde{\mathbf{x}} \geq 0} h(\tilde{\mathbf{w}}[1:d]^\top \mathbf{w}^*, \tilde{\mathbf{w}}[d+1]) \right] - \phi(\mathbf{w}^{*\top} \tilde{\mathbf{x}}[1:d]) \right| \leq \epsilon$$

where  $\tilde{\mathbf{w}} \sim \mathcal{N}(0, \mathbf{I}_{d+1})$  is an  $d$ -dimensional Gaussian.

We also have  $\mathbb{E}_{\tilde{\mathbf{w}} \in \mathcal{N}(0, \mathbf{I}_{d+1})} [h(\tilde{\mathbf{w}}[1:d]^\top \mathbf{w}^*, \tilde{\mathbf{w}}[d+1])^2] \leq (C_s(\phi, a))^2 (a^2 + 1)$ .

Proof: The statement can be found in Lemma B.1. from Li et al. (2022a)

## B Concept Class

To quantify complexity in (8), we define model complexity  $\mathcal{C}_\epsilon$  and sample complexity  $\mathcal{C}_s$  as in Li et al. (2022a) (Section 1.2) and Allen-Zhu & Li (2019) (Section 4). For a smooth function  $\phi(z) = \sum_{i=0}^{\infty} c_i z^i$ :

$$\mathcal{C}_\epsilon(\phi, R) = \sum_{i=0}^{\infty} \left( (C^* R)^i + \left( \frac{\sqrt{\log(1/\epsilon)}}{\sqrt{i}} C^* R \right)^i \right) |c_i|, \quad (19)$$

$$\mathcal{C}_s(\phi, R) = C^* \sum_{i=0}^{\infty} (i+1)^{1.75} R^i |c_i|, \quad (20)$$

where  $R \geq 0$  and  $C^*$  is a sufficiently large constant. These two quantities are used in the model complexity and sample complexity, which represent the required number of model parameters and training samples to learn  $\phi$  up to  $\epsilon$  error, respectively. Many population functions have bounded complexity. For instance, if  $\phi(z)$  is  $\exp(z)$ ,  $\sin(z)$ ,  $\cos(z)$  or polynomials of  $z$ , then  $\mathcal{C}_\epsilon(\phi, O(1)) \leq O(\text{poly}(1/\epsilon))$  and  $\mathcal{C}_s(\phi, O(1)) \leq O(1)$ .

Thus, the complexities of  $\mathcal{F}$  and  $\mathcal{G}$  are given by the tuples  $(p_{\mathcal{F}}, \mathcal{C}_s(\mathcal{F}), \mathcal{C}_\epsilon(\mathcal{F}))$  and  $(p_{\mathcal{G}}, \mathcal{C}_s(\mathcal{G}), \mathcal{C}_\epsilon(\mathcal{G}))$ .  $\mathcal{F}$  and  $\mathcal{G}$  are composed by  $p_{\mathcal{F}}$  and  $p_{\mathcal{G}}$  different smooth functions.

We also state some simple properties regarding our complexity measure. We define  $\mathfrak{B}_{\mathcal{F}} := \max_n \|\mathcal{F}_{n, A^*}(X)\|_2$ ,  $\mathfrak{B}_{\mathcal{F} \circ \mathcal{G}} := \max_n \|\mathcal{G}_{n, A^*}(\mathcal{F}_{A^*}(X))\|_2$  for all  $X$  satisfying  $\|x_n\| = 1$ . Assume  $\mathcal{G}(\cdot)$  is  $\mathcal{L}_{\mathcal{G}}$  Lipschitz continuous. It is simple to verify (see Lemma A.5) that  $\mathfrak{B}_{\mathcal{F}} \leq \sqrt{k} p_{\mathcal{F}} \mathcal{C}_s(\mathcal{F}, \|A^*\|_1) \|A^*\|_1$ ,  $\mathcal{L}_{\mathcal{G}} \leq \sqrt{k} p_{\mathcal{G}} \mathcal{C}_s(\mathcal{G}, \mathfrak{B}_{\mathcal{F}} \|A^*\|_1)$  and  $\mathfrak{B}_{\mathcal{F} \circ \mathcal{G}} \leq k p_{\mathcal{F}} \mathcal{C}_s(\mathcal{F}, \|A^*\|_1) \cdot p_{\mathcal{G}} \mathcal{C}_s(\mathcal{G}, \mathfrak{B}_{\mathcal{F}} \|A^*\|_1) \mathfrak{B}_{\mathcal{F}} \|A^*\|_1$ .

## C Theorem 3.2 Proof Details

Let us define learner networks that are single-skip two-hidden-layer with ReLU activation  $\text{out}_n : \mathbb{R}^{d \times N} \times \mathbb{R}^N \rightarrow \mathbb{R}^k$  with

$$\text{out}_n(X, A; W, V) = \text{out}_n^1(X, A) + C\sigma(V \text{out}_n^1(X, A) a_n) \quad (21)$$

and

$$\text{out}_n^1(X, A) = CD_W \odot (WXA^1) \in R^{k \times N}, \quad (22)$$

$$\text{out}_n^1(X, A) = CD_W^n WXA_n^1 \in R^k \quad (23)$$

where

- $A = [a_1, a_2, \dots, a_N] \in R^{N \times N}$  denotes the normalized adjacency matrix.
- $X \in R^{d \times N}$  denotes the the matrix of  $d$  dimension features of  $N$  nodes.
- $W \in R^{m \times d}$  denotes the first hidden layer weight.

- $V \in R^{m \times k}$  denotes the second hidden layer weight.
- $C \in R^{k \times m}$  denotes the output layer weight.
- $D_W = [\mathbf{1}_{(WX_{a_1} \geq 0)}, \mathbf{1}_{(WX_{a_2} \geq 0)}, \dots, \mathbf{1}_{(WX_{a_N} \geq 0)}]$ ,  $D_W^n = \text{diag}\{\mathbf{1}_{(WX_{a_n} \geq 0)}\}$
- $D_V = \text{diag}\{\mathbf{1}_{(V \text{out}_n^1(X,A)_{a_1} \geq 0)}, \mathbf{1}_{(V \text{out}_n^1(X,A)_{a_2} \geq 0)}, \dots, \mathbf{1}_{(V \text{out}_n^1(X,A)_{a_N} \geq 0)}\}$ ,  $D_V^n = \text{diag}\{\mathbf{1}_{(V \text{out}_n^1(X,A)_{a_n} \geq 0)}\}$

The  $l_2$  loss is represented as a function of the weight deviations  $\mathbf{W}, \mathbf{V}$  from initiation  $W^{(0)}$  and  $V^{(0)}$ , i.e.,

$$L(\mathbf{W}, \mathbf{V}) = \text{Obj}_n(X, A^{1t}, A^{2t}, y_n; W^{(0)} + \mathbf{W}, V^{(0)} + \mathbf{V}). \quad (24)$$

Let  $W^{(t)} = W^{(0)} + \mathbf{W}_t$ ,  $V^{(t)} = V^{(0)} + \mathbf{V}_t$ . We assume  $0 < \alpha \leq \tilde{O}\left(\frac{1}{k p_{\mathcal{G}} \mathfrak{C}_s(\mathcal{G}, \mathfrak{B}_{\mathcal{F}} \|A^*\|_1)}\right)$  throught the training. We prove that  $\|\mathbf{W}_t\|_2$  and  $\|\mathbf{V}_t\|_2$  are bounded by  $\tau_w$  and  $\tau_v$  in Table 3 during training, i.e.,  $\|\mathbf{W}_t\|_2 \leq \tau_w$ ,  $\|\mathbf{V}_t\|_2 \leq \tau_v$  for all  $t$ . See Appendix C.4 for the proof.

Table 3: Parameter choices

$\sigma_w$	$m^{-\frac{1}{2}+0.01} \leq \sigma_w \leq m^{-0.01}$	$\sigma_v$	$\sigma_v = \Theta(\text{polylog}(m))$
$\tau_w$	$\tau_w = \tilde{\Theta}(k p_{\mathcal{F}} \mathfrak{C}_s(\mathcal{F}, \ A^*\ _1))$ and $m^{\frac{1}{8}+0.001} \sigma_w \leq \tau_w \leq m^{\frac{1}{8}-0.001} \sigma_w^{\frac{1}{4}}$		
$\tau_v$	$\tau_v = \tilde{\Theta}(\alpha k p_{\mathcal{G}} \mathfrak{C}_s(\mathcal{G}, \mathfrak{B}_{\mathcal{F}} \ A^*\ _1) \mathfrak{B}_{\mathcal{F}} \ A^*\ _1)$ and $\sigma_v \left(\frac{k}{m}\right)^{\frac{3}{8}} \leq \tau_v \leq \frac{\sigma_v}{\text{polylog}(m) \ A^*\ _1} < 1$		

### C.1 Coupling

**Lemma C.1.** *We show that the weights after a properly bounded amount of updates stay close to the initialization, and many good properties occur. Suppose that  $\|\mathbf{W}\|_2 \leq \tau_w$ ,  $\|\mathbf{V}\|_2 \leq \tau_v$ ,  $W_0$  from  $\mathcal{N}(0, \sigma_w^2)$  and  $V_0$  from  $\mathcal{N}(0, \sigma_v^2/m)$ , we have that*

1. 
$$\|D_{\mathbf{W}^{(0)}}^n - D_{\mathbf{W}+\mathbf{W}_0}^n\|_0 \leq O\left(\left(\frac{\tau_w}{\sigma_w}\right)^{2/3} m^{2/3}\right) \quad (25)$$

2. 
$$\|CD_{\mathbf{W}+\mathbf{W}_0}^n \mathbf{W} X_{a_n} - CD_{\mathbf{W}+\mathbf{W}_0}^n (\mathbf{W} + \mathbf{W}_0) X_{a_n}\|_2 \leq \tilde{O}\left(\frac{\sqrt{s}}{\sqrt{m}} \tau_w \|A\|_1 + \sqrt{k} \sigma_w \|A\|_1\right) \quad (26)$$

3. 
$$\|\text{out}_n^1(X, A; W)\|_2 \leq \tilde{O}(\tau_w \|A\|_1) \quad (27)$$

4. 
$$\|D_{\mathbf{V}^{(0)}}^n - D_{\mathbf{V}+\mathbf{V}_0}^n\|_0 \leq O\left(\left(\frac{\tau_v}{\sigma_v}\right)^{2/3} m\right) \quad (28)$$

5. 
$$\begin{aligned} & \|CD_{\mathbf{V}+\mathbf{V}_0}^n \mathbf{V} \text{out}_n^1(X, A)_{a_n} - CD_{\mathbf{V}+\mathbf{V}_0}^n (\mathbf{V} + \mathbf{V}_0) \text{out}_n^1(X, A)_{a_n}\|_2 \\ & \leq \tilde{O}\left(\left(\frac{\sqrt{k}}{\sqrt{m}} \sigma_v + \frac{\sqrt{s}}{\sqrt{m}} \tau_v\right) \|\text{out}_n^1(X, A)\|_2 \|A\|_1\right) \end{aligned} \quad (29)$$

6. 
$$\|CD_{\mathbf{V}+\mathbf{V}_0}^n \mathbf{V}_0\|_2 \leq \tau_v \left(\frac{\tau_v}{\sigma_v}\right)^{1/3} \quad (30)$$

7.

$$\|CD_{\mathbf{V}+\mathbf{V}_0}^n(\mathbf{V}+\mathbf{V}_0)\text{out}^1(X,A)a_n\|_2 \leq \tilde{O}(\tau_v \|\text{out}_n^1(X,A)\|_2 \|A\|_1) \quad (31)$$

Proof:

1.  $\|WXa_n\|_2 \leq \|W\|_2 \|Xa_n\|_2 \leq \tau_w \|Xa_n\|_2$  and  $\langle W_0, Xa_n \rangle_j \sim \mathcal{N}(0, \|Xa_n\|_2^2 \sigma_w^2)$ , using Lemma A.2, we have

$$\|D_{\mathbf{W}^{(0)}}^n - D_{\mathbf{W}+\mathbf{W}_0}^n\|_0 \leq O\left(\left(\frac{\tau_w \|Xa_n\|_2}{\sigma_w \|Xa_n\|_2 \sqrt{m}}\right)^{2/3} m\right) \quad (32)$$

2. We write  $CD_{\mathbf{W}+\mathbf{W}_0}^n \mathbf{W}Xa_n^{*1} - CD_{\mathbf{W}+\mathbf{W}_0}^n (\mathbf{W}+\mathbf{W}_0)Xa_n^{*1} = -CD_{\mathbf{W}_0}^n \mathbf{W}^{(0)}Xa_n^{*1} + \mathbf{C}(D_{\mathbf{W}_0}^n - D_{\mathbf{W}+\mathbf{W}_0}^n) \mathbf{W}_0Xa_n^{*1}$ . For the first term,  $\|D_{\mathbf{W}_0}^n \mathbf{W}_0Xa_n\|_2 \leq \|\mathbf{W}_0Xa_n\|_2 \leq O(\sigma_w \|A\|_1 \sqrt{m})$ , so  $\|CD_{\mathbf{W}_0}^n \mathbf{W}_0Xa_n\|_2 \leq \tilde{O}(\sqrt{k}\sigma_w \|A\|_1)$

For the second term, using Lemma A.2 again, we have

$$\|(D_{\mathbf{W}_0}^n - D_{\mathbf{W}+\mathbf{W}_0}^n) \mathbf{W}_0Xa_n\|_2 \leq \|WXa_n\|_2 \leq \tau_w \|A\|_1$$

Using Lemma A.1, for every  $s$ -sparse vector  $\mathbf{y}$ , it satisfies  $\|\mathbf{A}\mathbf{y}\|_2 \leq e^{O(\sqrt{\frac{s}{m}})} \|\mathbf{y}\|_2$  with high probability. The sparsity of the second term is  $s = (\frac{\tau_w}{\sigma_w \|A\|_1})^{2/3} m^{2/3}$ , so we have

$$\|\mathbf{C}(D_{\mathbf{W}_0}^n - D_{\mathbf{W}+\mathbf{W}_0}^n) \mathbf{W}_0Xa_n\|_2 \leq \tilde{O}\left(\frac{\sqrt{s}}{\sqrt{m}}\right) \cdot \|WXa_n\|_2 \leq \tilde{O}\left(\frac{\sqrt{s}}{\sqrt{m}} \tau_w \|A\|_1\right).$$

3.  $\|\text{out}_n^1(X,A)\|_2 \leq \|CD_{\mathbf{W}+\mathbf{W}_0}^n \mathbf{W}Xa_n\|_2 + \|CD_{\mathbf{W}+\mathbf{W}_0}^n \mathbf{W}Xa_n - CD_{\mathbf{W}+\mathbf{W}_0}^n (\mathbf{W}+\mathbf{W}_0)Xa_n\|_2$ . Using  $\|\mathbf{C}\|_2 \leq 1$  with high probability, we have  $\|CD_{\mathbf{W}+\mathbf{W}_0}^n Xa_n\|_2 \leq \tilde{O}(\tau_w \|A\|_1)$
4. Similar to (25),  $\|V\text{out}^1(X,A)a_n\|_2 \leq \tau_v \|\text{out}^1(X,A)a_n\|_2$  and  $\langle V_0, \text{out}^1(X,A)a_n \rangle_j \sim \mathcal{N}(0, \|\text{out}^1(X,A)a_n\|_2^2 \sigma_v^2)$ , using Lemma A.2 we can prove it.
5. We write  $CD_{\mathbf{V}+\mathbf{V}_0}^n \mathbf{V}\text{out}^1(X,A)a_n - CD_{\mathbf{V}+\mathbf{V}_0}^n (\mathbf{V}+\mathbf{V}_0)\text{out}^1(X,A)a_n = -CD_{\mathbf{V}_0}^n \mathbf{V}^{(0)}\text{out}^1(X,A)a_n + \mathbf{C}(D_{\mathbf{V}_0}^n - D_{\mathbf{V}+\mathbf{V}_0}^n) \mathbf{V}_0\text{out}^1(X,A)a_n$ . Similar to (26), we have  $\|CD_{\mathbf{V}_0}^n \mathbf{V}^{(0)}\text{out}^1(X,A)a_n\|_2 \leq \tilde{O}(\sqrt{k}/\sqrt{m}) \cdot O(\sigma_v \|\text{out}^1(X,A)a_n\|_2)$  and  $\|\mathbf{C}(D_{\mathbf{V}_0}^n - D_{\mathbf{V}+\mathbf{V}_0}^n) \mathbf{V}_0\text{out}^1(X,A)a_n\|_2 \leq \tilde{O}\left(\frac{\sqrt{s}}{\sqrt{m}} \tau_v \|\text{out}^1(X,A)a_n\|_2\right)$ .  $\|\text{out}^1(X,A)a_n\|_2 \leq \|\text{out}_n^1(X,A)\|_2 \|A\|_1$
6. From 5, it is easy to get.
7. From 3, it is easy to get.

## C.2 Existential

Consider random function  $S_n((X,A); \mathbf{W}^*) = (S_n^1((X,A); \mathbf{W}^*), \dots, S_n^k((X,A); \mathbf{W}^*))$  in which

$$S_n^r((X,A); \mathbf{W}^*) \stackrel{\text{def}}{=} \sum_{i=1}^m a_{r,i} \cdot \langle w_i^*, Xa_n \rangle \cdot \mathbf{1}_{\langle w_i^{(0)}, Xa_n \rangle \geq 0} \quad (33)$$

where  $W^*$  is a given matrix,  $W^0$  is a random matrix where each  $w_i^{(0)}$  is i.i.d. from  $\mathcal{N}(0, \frac{\mathbf{I}}{m})$  and  $a_{r,i}$  is i.i.d. from  $\mathcal{N}(0, 1)$ .

Based on Lemma B.1. from Li et al. (2022a) and Lemma E.1. from ?, we have

**Lemma C.2.** *Given any  $\mathcal{F} : \mathbb{R}^d \rightarrow \mathbb{R}^k$  with general complexity  $(p, \mathfrak{C}_s(\mathcal{F}, \|A\|_1) \|A\|_1, \mathfrak{C}_\varepsilon(\mathcal{F}, \|A\|_1) \|A\|_1)$ , for every  $\varepsilon \in \left(0, \frac{1}{pk\mathfrak{C}_s(\mathcal{F}, \|A\|_1) \|A\|_1}\right)$ , there exist  $M = \text{poly}(\mathfrak{C}_\varepsilon(\mathcal{F}, \|A\|_1), \|A\|_1, 1/\varepsilon)$  such that if  $m \geq M$ , then with high probability there is a construction  $\mathbf{W}^* = (w_1^*, \dots, w_m^*) \in \mathbb{R}^{m \times d}$  with*

$$\|\mathbf{W}^*\|_{2,\infty} \leq \frac{kp\mathfrak{C}_\varepsilon(\mathcal{F}, \|A\|_1) \|A\|_1}{m} \quad \text{and} \quad \|\mathbf{W}^*\|_F \leq \tilde{O}\left(\frac{kp\mathfrak{C}_s(\mathcal{F}, \|A\|_1) \|A\|_1}{\sqrt{m}}\right) \quad (34)$$

satisfying, for every  $x_n \in \mathbb{R}^d$  and  $\|x_n\|_2 \leq 1$ , with probability at least  $1 - e^{-\Omega(\sqrt{m})}$

$$\sum_{r=1}^k |\mathcal{F}_n^r(X, A) - S_n^r(X, A; W^*)| \leq \varepsilon \cdot \|A\|_1 \quad (35)$$

where  $G_n(X, A; W^*) = \begin{bmatrix} S_n^1(X, A; W^*) \\ \vdots \\ S_n^k(X, A; W^*) \end{bmatrix}$  and  $S_n(X, A; W^*) = CD_{W+W_0}^n W^* X a_n$

Proof: Define  $w_j^* = \sum_{r \in [k]} a_{r,j} \sum_{i \in [p]} a_{r,i}^* h^{(r,i)} \left( \sqrt{m} \langle w_j^{(0)}, w_{1,i}^* \rangle \right)$   $w_{2,i}^*$  has the same distribution with  $\alpha_1$  in Lemma A.6.

Using Lemma A.6 we have  $|h^{(r,i)}| \leq \mathfrak{C}_\varepsilon(\mathcal{F}, \|A\|_1) \|A\|_1$  and using Lemma E.1. from ?, we have for our parameter choice of  $m$ , with probability at least  $1 - e^{-\Omega(m\varepsilon^2/(k^4 p^2 \mathfrak{C}_\varepsilon(\mathcal{F}, \|A\|_1) \|A\|_1))}$

$$|\mathcal{F}_n^r(X, A) - S_n^r(X, A; W^*)| \leq \frac{\varepsilon}{k}.$$

We have for each  $j \in [m]$ , with high probability  $\|w_j^*\|_2 \leq \tilde{O}\left(\frac{kp\mathfrak{C}_\varepsilon(\mathcal{F}, \|A\|_1) \|A\|_1}{m}\right)$ . This means  $\|\mathbf{W}^*\|_{2,\infty} \leq \tilde{O}\left(\frac{kp\mathfrak{C}_\varepsilon(\mathcal{F}, \|A\|_1) \|A\|_1}{m}\right)$ . As for the Frobenius norm,

$$\|\mathbf{W}^*\|_F^2 = \sum_{j \in [m]} \|w_j^*\|_2^2 \leq \sum_{j \in [m]} \tilde{O}\left(\frac{k^2 p}{m^2}\right) \cdot \sum_{i \in [p]} h^{(r,i)} \left( \sqrt{m} \langle w_j^{(0)}, w_{1,i}^* \rangle \right)^2 \quad (36)$$

Applying Hoeffding's concentration, we have with probability at least  $1 - e^{-\Omega(\sqrt{m})}$

$$\begin{aligned} \sum_{j \in [m]} h^{(r,i)} \left( \sqrt{m} \langle w_j^{(0)}, w_{1,i}^* \rangle, \sqrt{m} b_j^{(0)} \right)^2 &\leq m \cdot (\mathfrak{C}_s(\mathcal{F}, \|A\|_1) \|A\|_1^2), \\ &+ m^{3/4} \cdot (\mathfrak{C}_\varepsilon(\mathcal{F}, \|A\|_1) \|A\|_1)^2, \\ &\leq 2m (\mathfrak{C}_s(\mathcal{F}, \|A\|_1) \|A\|_1)^2. \end{aligned} \quad (37)$$

Putting this back to (36) we have  $\|\mathbf{W}^*\|_F^2 \leq \tilde{O}\left(\frac{k^2 p^2 (\mathfrak{C}_s(\mathcal{F}, \|A\|_1) \|A\|_1)^2}{m}\right)$ .

**Lemma C.3.** *Under the assumptions of Lemma C.1, suppose  $\alpha \in (0, 1)$  and  $\tilde{\alpha} = \frac{\alpha}{k(p_{\mathcal{F}} \mathfrak{C}_s(\mathcal{F}, \|A\|_1) + p_{\mathcal{G}} \mathfrak{C}_s(\mathcal{G}, \|A\|_1))}$ , there exist  $M = \text{poly}(\mathfrak{C}_\alpha(\mathcal{F}, \|A\|_1), \mathfrak{C}_\alpha(\mathcal{G}, \|A\|_1), \|A\|_1, \tilde{\alpha}^{-1})$  satisfying that for every  $m \geq M$ ,  $\|\mathbf{W}^*\|_F \leq \tilde{O}(kp_{\mathcal{F}} \mathfrak{C}_s(\mathcal{F}))$  and  $\|\mathbf{V}^*\|_F \leq \tilde{O}(\tilde{\alpha} kp_{\mathcal{G}} \mathfrak{C}_s(\mathcal{G}))$  with high probability*

1. 
$$\mathbb{E}_{n \in \mathcal{V}, (X, y_n) \sim \mathcal{D}} [\|CD_{\mathbf{W}_0}^n \mathbf{W}^* X a_n - \mathcal{F}_n(X, A)\|_2] \leq \tilde{\alpha}^2 \|A\|_1 \quad (38)$$

2. 
$$\|CD_{\mathbf{V}_0}^n \mathbf{V}^* \text{out}^1(X, A) a_n - \alpha \mathcal{G}_n(\text{out}^1(X, A), A)\|_2 \leq \tilde{\alpha}^2 \cdot \|\text{out}_n^1(X, A)\|_2 \|A\|_1 \quad (39)$$

3. 
$$\mathbb{E}_{n \in \mathcal{V}, (X, y_n) \sim \mathcal{D}} [\|CD_{\mathbf{W}}^n \mathbf{W}^* X a_n - \mathcal{F}_n(X, A)\|_2] \leq O(\tilde{\alpha}^2 \|A\|_1) \quad (40)$$

4. 
$$\begin{aligned} &\|CD_{\mathbf{V}}^n \mathbf{V}^* \text{out}^1(X, A) a_n - \alpha \mathcal{G}_n(\text{out}^1(X, A), A)\|_2 \\ &\leq \left( \tilde{\alpha}^2 + O\left(\tau_v \left(\frac{\tau_v}{\sigma_v}\right)^{1/3}\right) \right) \|\text{out}_n^1(X, A)\|_2 \|A\|_1 \end{aligned} \quad (41)$$

5.

$$\mathbb{E}_{n \in \mathcal{V}, (X, y_n) \sim \mathcal{D}} [\|CD_{\mathbf{W}_0 + \mathbf{W}}^n (\mathbf{W}^* - \mathbf{W}) X a_n - (\mathcal{F}_n(X, A) - \text{out}_n^1(X, A))\|_2] \leq \tilde{\alpha}^2 \|A\|_1 \quad (42)$$

Proof:

1. Using Lemma C.2, we can find a  $\mathbf{W}^*$  satisfying  $\|CD_{\mathbf{W}_0 + \mathbf{W}}^n \mathbf{W}^* X a_n - \mathcal{F}_n(X, A)\|_2$  small enough with probability at least  $1 - e^{-\Omega(\sqrt{m})}$ .
2. Using Lemma C.2 and  $\|\text{out}^1(X, A) a_n\|_2 \leq \|\text{out}_n^1(X, A)\|_2 \|A\|_1$ , we can easily prove it.
3.  $\|\mathbf{W}^* X a_n\|_2 \leq O(\|\mathbf{W}^*\|_F \|X a_n\|_2) \leq O(\tau_w \|A\|_1)$ .  $\|C(D_{\mathbf{W}}^n - D_{\mathbf{W}_0}^n) \mathbf{W}^* X a_n\|_2 \leq O(\sqrt{s} \tau_w \|A\|_1 / \sqrt{m})$  where  $s$  is the maximum sparsity of  $(D_{\mathbf{W}}^n - D_{\mathbf{W}_0}^n)$ . Using (25), we know  $s \leq O\left(\left(\frac{\tau_w}{\sigma_w}\right)^{2/3} m^{2/3}\right)$ . This, combining with (38) gives

$$\begin{aligned} \mathbb{E}_{n \in \mathcal{V}, (X, y_n) \sim \mathcal{D}} [\|CD_{\mathbf{W}}^n \mathbf{W}^* X a_n - \mathcal{F}_n(X, A)\|_2] &\leq \tilde{\alpha}^2 \|A\|_1 + O\left(\tau_w \left(\frac{\tau_w}{\sigma_w}\right)^{1/3} / m^{1/6}\right) \\ &\leq O(\tilde{\alpha}^2 \|A\|_1) \end{aligned} \quad (43)$$

4. Using (28) and  $\|\mathbf{V}^* \text{out}^1(X, A) a_n\|_2 \leq O(\tau_v \|\text{out}_n^1(X, A)\|_2 \|A\|_1)$  we can easily prove it.
5. Using (26) and (40), with larger enough  $m$ , we can prove it.

### C.3 Optimization

We write the gradient of loss function as  $\nabla_{\mathbf{W}} \text{Obj}_n(\mathbf{W}) = \nabla_{\mathbf{W}} \text{Obj}_n^1(\mathbf{W}) + \nabla_{\mathbf{W}} \text{Obj}_n^2(\mathbf{W})$ , where  $\nabla_{\mathbf{W}} \text{Obj}_n^1(\mathbf{W}) = \nabla_{\mathbf{W}} \text{out}_n^1(X, A)$  and  $\nabla_{\mathbf{W}} \text{Obj}_n^2(\mathbf{W}) = \nabla_{\mathbf{W}} CD_{\mathbf{V}}^n V \text{out}^1(X, A) a_n$ , we can write its gradient as follows.

$$\begin{aligned} \langle \nabla_{\mathbf{W}} \text{Obj}_n^1(\mathbf{W}), -\mathbf{W}' \rangle &= \text{tr}(X a_n (y_n - \text{out}_n(X, A))^\top CD_{\mathbf{W} + \mathbf{W}_0}^n \mathbf{W}') \\ &= \text{tr}((y_n - \text{out}_n(X, A))^\top CD_{\mathbf{W} + \mathbf{W}_0}^n \mathbf{W}' X a_n) \\ &= \langle y_n - \text{out}_n(X, A), CD_{\mathbf{W} + \mathbf{W}_0}^n \mathbf{W}' X a_n \rangle \end{aligned} \quad (44)$$

$$\begin{aligned} \langle \nabla_{\mathbf{W}} \text{Obj}_n^2(\mathbf{W}), -\mathbf{W}' \rangle &= \text{tr}\left(\sum_{i=1}^N a_{ni} X a_i (y_n - \text{out}_n(X, A))^\top CD_{\mathbf{V} + \mathbf{V}_0}^n (\mathbf{V}^{(0)} + \mathbf{V}) CD_{\mathbf{W} + \mathbf{W}_0}^i \mathbf{W}'\right) \\ &= \text{tr}\left(\sum_{i=1}^N a_{ni} (y_n - \text{out}_n(X, A))^\top CD_{\mathbf{V} + \mathbf{V}_0}^n (\mathbf{V}^{(0)} + \mathbf{V}) CD_{\mathbf{W} + \mathbf{W}_0}^i \mathbf{W}' X a_i\right) \\ &= \text{tr}((y_n - \text{out}_n(X, A))^\top \sum_{i=1}^N a_{ni} CD_{\mathbf{V} + \mathbf{V}_0}^n (\mathbf{V}^{(0)} + \mathbf{V}) CD_{\mathbf{W} + \mathbf{W}_0}^i \mathbf{W}' X a_i) \\ &= \langle y_n - \text{out}_n(X, A), CD_{\mathbf{V} + \mathbf{V}_0}^n (\mathbf{V}^{(0)} + \mathbf{V}) C(D_{\mathbf{W} + \mathbf{W}_0} \odot \mathbf{W}' X A) a_n \rangle \end{aligned} \quad (45)$$

$$\begin{aligned} \langle \nabla_{\mathbf{V}} \text{Obj}_n(\mathbf{V}), -\mathbf{V}' \rangle &= \text{tr}(\text{out}(X) a_n (y_n - \text{out}_n(X, A))^\top CD_{\mathbf{V} + \mathbf{V}_0}^n \mathbf{V}') \\ &= \text{tr}((y_n - \text{out}_n(X, A))^\top CD_{\mathbf{V} + \mathbf{V}_0}^n \mathbf{V}' \text{out}(X) a_n) \\ &= \langle y_n - \text{out}_n(X, A), CD_{\mathbf{V} + \mathbf{V}_0}^n \mathbf{V}' \text{out}(X) a_n \rangle \end{aligned} \quad (46)$$

Let us define  $f(\mathbf{W}') = CD_{\mathbf{W} + \mathbf{W}_0}^n \mathbf{W}' X a_n + CD_{\mathbf{V} + \mathbf{V}_0}^n (\mathbf{V}^{(0)} + \mathbf{V}) C(D_{\mathbf{W} + \mathbf{W}_0} \odot \mathbf{W}' X A) a_n$  and  $g(\mathbf{V}') = CD_{\mathbf{V} + \mathbf{V}_0}^n \mathbf{V}' \text{out}(X) a_n$ . Therefore,

$$\langle \nabla_{\mathbf{W}, \mathbf{V}} \text{Obj}_n(\mathbf{W}, \mathbf{V}), (-\mathbf{W}', -\mathbf{V}') \rangle = \langle y_n - \text{out}_n(X, A), f(\mathbf{W}') + g(\mathbf{V}') \rangle \quad (47)$$

**Claim C.4.** We have that for all  $\mathbf{W}$  and  $\mathbf{V}$  satisfying  $\|\mathbf{W}\|_F \leq \tau_w$  and  $\|\mathbf{V}\|_F \leq \tau_v$ , it holds that

$$\|\nabla_{\mathbf{W}} \text{Obj}(\mathbf{W}, \mathbf{V}; (x, y))\|_F \leq \|A\|_1 \|y_n - \text{out}_n(X, A)\|_2 \cdot O(\sigma_v + 1) \quad (48)$$

$$\|\nabla_{\mathbf{V}} \text{Obj}(\mathbf{W}, \mathbf{V}; (x, y))\|_F \leq \tau_w \|A\|_1 \|y_n - \text{out}_n(X, A)\|_2 \cdot O(1) \quad (49)$$

Proof:

$$\begin{aligned} \|\nabla_{\mathbf{W}} \text{Obj}(\mathbf{W}, \mathbf{V}; (x, y))\|_F &= \|X a_n (y_n - \text{out}_n(X, A))^\top \\ &\quad \times (CD_{W+W_0}^n + CD_{V+V_0}^n(\mathbf{V}^{(0)} + \mathbf{V})CD_{W+W_0}^n)\|_F \\ &\leq \|X a_n\|_2 \|y_n - \text{out}_n(X, A)\|_2 \\ &\quad \times \|CD_{W+W_0}^n + CD_{V+V_0}^n(\mathbf{V}^{(0)} + \mathbf{V})CD_{W+W_0}^n\|_2 \\ &\leq \|A\|_1 \|y_n - \text{out}_n(X, A)\|_2 \cdot O(\sigma_v + 1) \end{aligned} \quad (50)$$

In (50), the last inequality uses  $\|V^{(0)}\|_2 = O(\tau_v)$  and  $\|C\|_2 \leq 1$ .

$$\begin{aligned} \|\nabla_{\mathbf{V}} \text{Obj}(\mathbf{W}, \mathbf{V}; (x, y))\|_F &= \|\text{out}_n(X, A) a_n (y_n - \text{out}_n(X, A))^\top CD_{V+V_0}^n\|_F \\ &\leq \|\text{out}_n(X, A) a_n\|_2 \|y_n - \text{out}_n(X, A)\|_2 \|CD_{V+V_0}^n\|_2 \\ &\leq \tau_w \|A\|_1 \|y_n - \text{out}_n(X, A)\|_2 \cdot O(1) \end{aligned} \quad (51)$$

In (51), the last inequality uses (27) and  $\|C\|_2 \leq 1$ .

**Claim C.5.** In the setting of Lemma C.1, we have  $f(\mathbf{W}^* - \mathbf{W}) + g(\mathbf{V}^* - \mathbf{V}) = \mathcal{H}_{n, A^*}(X, A) - \text{out}_n(X, A) + \text{Err}_n$  with

$$\begin{aligned} \mathbb{E}_{n \in \mathcal{V}, (X, y_n) \sim \mathcal{D}} \|\text{Err}_n\|_2 &\leq \mathbb{E}_{n \in \mathcal{V}, (X, y_n) \sim \mathcal{D}} [O(\tau_v \|A\|_1 + \alpha \mathfrak{L}_{\mathcal{G}} \|A\|_1) \cdot \|\mathcal{H}_{n, A^*}(X, A) - \text{out}_n(X, A)\|_2] \\ &\quad + O\left(\tau_v^2 \|A\|_1^2 \mathfrak{B}_{\mathcal{F}} + \tau_v \tilde{\alpha}^2 \|A\|_1^2 + \alpha \tau_v \|A\|_1^2 \mathfrak{L}_{\mathcal{G}} \mathfrak{B}_{\mathcal{F}}\right) \end{aligned} \quad (52)$$

Proof: Based on the definition of  $f(\mathbf{W}')$  and  $g(\mathbf{V}')$ , we have

$$\begin{aligned} f(\mathbf{W}^* - \mathbf{W}; X, a_n^*) + g(\mathbf{V}^* - \mathbf{V}; X, a_n^*) &= CD_{W+W_0}^n (W^* - W) X a_n \\ &\quad + CD_{V+V_0}^n (V^* - V) \text{out}(X) a_n \\ &\quad + CD_{V+V_0}^n (\mathbf{V}^{(0)} + \mathbf{V}) C(D_{W+W_0} \odot (W^* - W) X A) a_n \\ &= \underbrace{CD_{V+V_0}^n (\mathbf{V}^{(0)} + \mathbf{V}) C(D_{W+W_0} \odot (W^* - W) X A) a_n}_{\clubsuit} \\ &\quad + \underbrace{CD_{W+W_0}^n W^* X a_n + CD_{V+V_0}^n V^* \text{out}^1(X, A) a_n}_{\spadesuit} \\ &\quad - \underbrace{(CD_{W+W_0}^n W X a_n + CD_{V+V_0}^n V \text{out}^1(X, A) a_n)}_{\diamond} \end{aligned} \quad (53)$$

1. For the  $\clubsuit$  term,

$$\begin{aligned} \clubsuit &\leq \left( \|CD_{V+V_0}^n \mathbf{V}^{(0)}\|_2 + \|C\|_2^2 \|\mathbf{V}\|_2 \right) \|C(D_{W+W_0} \odot (W^* - W) X A) a_n\|_2 \\ &\leq O(1) \cdot O(\tau_v) \cdot \sum_{i=1}^N a_{ni} (\|\mathcal{F}(x) - \text{out}_n^1(X, A)\|_2 + O(\tilde{\alpha}^2 \|A\|_1)) \\ &\leq O(\tau_v) \left( \|\mathcal{F}(x) - \text{out}_i(x)\|_2 \|A\|_1 + O(\tilde{\alpha}^2 \|A\|_1^2) \right) \end{aligned} \quad (54)$$

together with  $\tau_v \leq \frac{1}{\text{polylog}(m)} \sigma_v$ .

2. For the  $\spadesuit$  term,

$$\begin{aligned} \spadesuit - (\mathcal{F}_n(X, A) + \alpha\mathcal{G}(\mathcal{F}(x), a_n)) &= CD_{W+W_0}^n W^* X a_n - \mathcal{F}_n(X, A) \\ &\quad + CD_{V+V_0}^n V^* \text{out}_n^1(X, A) a_n - \alpha\mathcal{G}(\text{out}_n^1(X, A), a_n) \\ &\quad + \alpha\mathcal{G}(\text{out}_n^1(X, A), a_n) - \alpha\mathcal{G}(\mathcal{F}(x), a_n) \end{aligned} \quad (55)$$

The first term uses (38), the second term uses (39) and the third term uses the Lipschitz continuity of  $\mathcal{G}$ , so we have

$$\begin{aligned} \|\spadesuit - (\mathcal{F}(x) + \alpha\mathcal{G}(\mathcal{F}(x)))\|_2 &\leq O\left(\tilde{\alpha}^2 + \tau_v \left(\frac{\tau_v}{\sigma_v}\right)^{1/3}\right) \cdot \|\text{out}_n^1(X, a_n)\|_2 \|A\|_1 \\ &\quad + O(\alpha\mathfrak{L}_{\mathcal{G}}) \|\mathcal{F}(X) a_n - \text{out}_n^1(X) a_n\|_2 \\ &\leq O(\tau_v^2) \cdot \|\text{out}_n^1(x)\|_2 \|A\|_1 + O(\alpha\mathfrak{L}_{\mathcal{G}}) \|\mathcal{F}_n(x) - \text{out}_n^1(x)\|_2 \|A\|_1 \end{aligned} \quad (56)$$

We use  $\frac{1}{\sigma_v} \leq \tau_v^2$  and definition of  $\tilde{\alpha}$ .

3. For the  $\diamond$  term,

$$\|\diamond - \text{out}_n(X)\|_2 \leq O\left(\|\text{out}_n^1(x)\|_2 \|A\|_1\right) \tau_v^2 \quad (57)$$

where the inequality uses (26) and (29).

In sum, we have

$$Err \stackrel{\text{def}}{=} f(\mathbf{W}^* - \mathbf{W}; x) + g(\mathbf{V}^* - \mathbf{V}; x) - (\mathcal{F}(x) + \alpha\mathcal{G}(\mathcal{F}(x)) - \text{out}_n(X, A)) \quad (58)$$

satisfy

$$\begin{aligned} \mathbb{E}_{n \in \mathcal{V}, (X, y_n) \sim \mathcal{D}} \|Err\|_2 &\leq \mathbb{E}_{n \in \mathcal{V}, (X, y_n) \sim \mathcal{D}} \left[ O(\tau_v \|A\|_1 + \alpha\mathfrak{L}_{\mathcal{G}} \|A\|_1) \right. \\ &\quad \left. \times \|\mathcal{F}(x) - \text{out}_n^1(X, A)\|_2 + O(\|\text{out}_n^1(x)\|_2 \|A\|_1) \tau_v^2 \right] \\ &\quad + O\left(\tau_v \tilde{\alpha}^2 \|A\|_1^2\right) \end{aligned} \quad (59)$$

Using  $\|\text{out}_n^1(x)\|_2 \leq \|\text{out}_n^1(X, A) - \mathcal{F}(x)\|_2 + \mathfrak{B}_{\mathcal{F}}$ , we have

$$\begin{aligned} \mathbb{E}_{n \in \mathcal{V}, (X, y_n) \sim \mathcal{D}} \|Err\|_2 &\leq \mathbb{E}_{n \in \mathcal{V}, (X, y_n) \sim \mathcal{D}} [O(\tau_v \|A\|_1 + \alpha\mathfrak{L}_{\mathcal{G}} \|A\|_1) \cdot \|\mathcal{H}(x) - \text{out}_n(X, A)\|_2] \\ &\quad + O\left(\tau_v^2 \|A\|_1 \mathfrak{B}_{\mathcal{F}} + \tau_v \tilde{\alpha}^2 \|A\|_1^2\right) \\ &\quad + O(\tau_v \|A\|_1 + \alpha\mathfrak{L}_{\mathcal{G}} \|A\|_1) \cdot (\tau_v \|A\|_1 \mathfrak{B}_{\mathcal{F}} + \alpha\mathfrak{B}_{\mathcal{F} \circ \mathcal{G}}) \end{aligned} \quad (60)$$

Using  $\mathfrak{B}_{\mathcal{F} \circ \mathcal{G}} \leq \sqrt{k} p_{\mathcal{G}} \mathfrak{C}_{\mathfrak{s}}(\mathcal{G}, \|A\|_1 \mathfrak{B}_{\mathcal{F}}) (\|A\|_1 \mathfrak{B}_{\mathcal{F}})^2 \mathfrak{B}_{\mathcal{F}} \leq \frac{\tau_v}{\alpha} \mathfrak{B}_{\mathcal{F}}$ , so we have

$$\begin{aligned} \mathbb{E}_{n \in \mathcal{V}, (X, y_n) \sim \mathcal{D}} \|Err\|_2 &\leq \mathbb{E}_{n \in \mathcal{V}, (X, y_n) \sim \mathcal{D}} [O(\tau_v \|A\|_1 + \alpha\mathfrak{L}_{\mathcal{G}} \|A\|_1) \cdot \|\mathcal{H}(x) - \text{out}_n(X, A)\|_2] \\ &\quad + O\left(\tau_v^2 \|A\|_1^2 \mathfrak{B}_{\mathcal{F}} + \tau_v \tilde{\alpha}^2 \|A\|_1^2 + \alpha\tau_v \|A\|_1^2 \mathfrak{L}_{\mathcal{G}} \mathfrak{B}_{\mathcal{F}}\right) \end{aligned} \quad (61)$$

**Claim C.6.** *In the setting of Lemma C.1, if  $\tau_v \|A\|_1 \leq \frac{1}{\text{polylog}(m)}$ , we have*

$$\|\text{out}_n^1(X, A) - \mathcal{F}_n(X)\|_2 \leq 2 \|\text{out}_n^1(X, A) - \mathcal{H}_n(X, A)\|_2 + \alpha\mathfrak{B}_{\mathcal{F} \circ \mathcal{G}} + \tilde{O}(\tau_v \|A\|_1 \mathfrak{B}_{\mathcal{F}}) \quad (62)$$

Proof: Using 31 and  $\|\mathcal{G}(\mathcal{F}(X), a_n)\|_2 \leq \mathfrak{B}_{\mathcal{F} \circ \mathcal{G}}$ , we have

$$\begin{aligned} \|\text{out}_n^1(X, A) - \mathcal{F}_n(X)\|_2 &\leq \|\text{out}_n^1(X, A) - \mathcal{H}_n(X, A)\|_2 + \alpha\mathfrak{B}_{\mathcal{F} \circ \mathcal{G}} \\ &\quad + \tilde{O}(\tau_v \|A\|_1 (\|\text{out}_n^1(X, a_n) - \mathcal{F}_n(X, A)\|_2 + \mathfrak{B}_{\mathcal{F}})) \end{aligned} \quad (63)$$

Using  $\tau_v \|A\|_1$  small enough, we have

$$\|\text{out}_n^1(X, A) - \mathcal{F}_n(X)\|_2 \leq 2 \|\text{out}_n^1(X, A) - \mathcal{H}_n(X, A)\|_2 + \alpha\mathfrak{B}_{\mathcal{F} \circ \mathcal{G}} + \tilde{O}(\tau_v \|A\|_1 \mathfrak{B}_{\mathcal{F}}) \quad (64)$$

#### C.4 Proof of Theorem 3.2

Proof: Using (47) and Claim C.5, in iteration  $t$ , we have

$$\begin{aligned} & \langle \nabla_{\mathbf{W}, \mathbf{V}} \text{Obj}_n(\mathbf{W}_t, \mathbf{V}_t), (\mathbf{W}_t - \mathbf{W}^*, \mathbf{V}_t - \mathbf{V}^*) \rangle \\ &= \langle y_n - \text{out}(\mathbf{W}_t, \mathbf{V}_t), f(\mathbf{W}^* - \mathbf{W}) + g(\mathbf{V}^* - \mathbf{V}) \rangle \\ &= \langle y_n - \text{out}_n(\mathbf{W}_t, \mathbf{V}_t), \mathcal{H}_{n, A^*}(X) - \text{out}_n(\mathbf{W}_t, \mathbf{V}_t) + \text{Err}_t \rangle \end{aligned} \quad (65)$$

We also have

$$\begin{aligned} \|\mathbf{W}_{t+1} - \mathbf{W}^*\|_F^2 &= \|\mathbf{W}_t - \eta_w \nabla_{\mathbf{W}} \text{Obj}_n(\mathbf{W}_t) - \mathbf{W}^*\|_F^2 \\ &= \|\mathbf{W}_t - \mathbf{W}^*\|_F^2 - 2\eta_w \langle \nabla_{\mathbf{W}} \text{Obj}_n(\mathbf{W}_t), \mathbf{W}_t - \mathbf{W}^* \rangle \\ &\quad + \eta_w^2 \|\nabla_{\mathbf{W}} \text{Obj}_n(\mathbf{W}_t)\|_F^2, \end{aligned} \quad (66)$$

$$\begin{aligned} \|\mathbf{V}_{t+1} - \mathbf{V}^*\|_F^2 &= \|\mathbf{V}_t - \eta_v \nabla_{\mathbf{V}} \text{Obj}_n(\mathbf{V}_t) - \mathbf{V}^*\|_F^2 \\ &= \|\mathbf{V}_t - \mathbf{V}^*\|_F^2 - 2\eta_v \langle \nabla_{\mathbf{V}} \text{Obj}_n(\mathbf{V}_t), \mathbf{V}_t - \mathbf{V}^* \rangle \\ &\quad + \eta_v^2 \|\nabla_{\mathbf{V}} \text{Obj}_n(\mathbf{V}_t)\|_F^2 \end{aligned} \quad (67)$$

Using Algorithm 1, we have  $\mathbf{W}_{t+1} = \mathbf{W}_t - \eta_w \nabla_{\mathbf{W}} \text{Obj}_n(\mathbf{W}_t, \mathbf{V}_t)$  and  $\mathbf{V}_{t+1} = \mathbf{V}_t - \eta_v \nabla_{\mathbf{V}} \text{Obj}_n(\mathbf{W}_t, \mathbf{V}_t)$ , so we have

$$\begin{aligned} & \langle \nabla_{\mathbf{W}, \mathbf{V}} \text{Obj}_t(\mathbf{W}_t, \mathbf{V}_t), (\mathbf{W} - \mathbf{W}^*, \mathbf{V} - \mathbf{V}^*) \rangle \\ &= \underbrace{\frac{\eta_w}{2} \|\nabla_{\mathbf{W}} \text{Obj}_t(\mathbf{W}_t, \mathbf{V}_t)\|_F^2 + \frac{\eta_v}{2} \|\nabla_{\mathbf{V}} \text{Obj}_t(\mathbf{W}_t, \mathbf{V}_t)\|_F^2}_{\heartsuit} \\ &\quad + \frac{1}{2\eta_w} \|\mathbf{W}_t - \mathbf{W}^*\|_F^2 - \frac{1}{2\eta_w} \|\mathbf{W}_{t+1} - \mathbf{W}^*\|_F^2 + \frac{1}{2\eta_v} \|\mathbf{V}_t - \mathbf{V}^*\|_F^2 - \frac{1}{2\eta_v} \|\mathbf{V}_{t+1} - \mathbf{V}^*\|_F^2 \end{aligned} \quad (68)$$

Using Claim C.4 and change all  $A$  to  $A^*$ , we have

$$\begin{aligned} \heartsuit &\leq O(\eta_w \sigma_v^2 + \eta_v \tau_w^2) \cdot \|A\|_1^2 \|y_n - \text{out}_n(X, A^*)\|_2^2 \\ &\leq O(\eta_w \sigma_v^2 + \eta_v \tau_w^2) \cdot \|A\|_1^2 \left( \|\mathcal{H}_{n, A^*}(X) - \text{out}_n(X, A^*)\|_2^2 + \|\mathcal{H}_{n, A^*}(X) - y_n\|_2^2 \right) \end{aligned} \quad (69)$$

Therefore, as long as  $O(\eta_w \sigma_v^2 + \eta_v \tau_w^2) \leq 0.1$ , it satisfies

$$\begin{aligned} \frac{1}{4} \|\mathcal{H}_{n, A^*}(X) - \text{out}_n(X, A^*)\|_2^2 &\leq 2 \|\text{Err}_t\|_2^2 + 4 \|\mathcal{H}_{n, A^*}(X) - y_n\|_2^2 \\ &\quad + \frac{1}{2\eta_w} \|\mathbf{W}_t - \mathbf{W}^*\|_F^2 - \frac{1}{2\eta_w} \|\mathbf{W}_{t+1} - \mathbf{W}^*\|_F^2 \\ &\quad + \frac{1}{2\eta_v} \|\mathbf{V}_t - \mathbf{V}^*\|_F^2 - \frac{1}{2\eta_v} \|\mathbf{V}_{t+1} - \mathbf{V}^*\|_F^2 \end{aligned} \quad (70)$$

After telescoping for  $t = 0, 1, \dots, T_0 - 1$ ,

$$\begin{aligned} & \frac{\|\mathbf{W}_{T_0} - \mathbf{W}^*\|_F^2}{2\eta_w T_0} + \frac{\|\mathbf{V}_{T_0} - \mathbf{V}^*\|_F^2}{2\eta_v T_0} + \frac{1}{2T_0} \sum_{t=0}^{T_0-1} \|\mathcal{H}_{n, A^*}(X) - \text{out}_n(X, A^*)\|_2^2 \\ &\leq \frac{\|\mathbf{W}^*\|_F^2}{2\eta_w T_0} + \frac{\|\mathbf{V}^*\|_F^2}{2\eta_v T_0} + \frac{O(1)}{T_0} \sum_{t=0}^{T_0-1} \|\text{Err}_t\|_2^2 + \|\mathcal{H}_{n, A^*}(X) - y_t\|_2^2. \end{aligned} \quad (71)$$

Using  $O(\tau_v \|A\|_1 + \alpha \mathcal{L}_G) \leq 0.1$ , we have

$$\frac{1}{4T} \sum_{t=0}^{T-1} \mathbb{E}_{n \in \mathcal{V}, (X, y_n) \sim \mathcal{Z}} \|\mathcal{H}_{n, A^*}(X) - \text{out}_n(X, A^*)\|_2^2 \leq \frac{\|\mathbf{W}^*\|_F^2}{2\eta_w T} + \frac{\|\mathbf{V}^*\|_F^2}{2\eta_v T} + O(\text{OPT} + \epsilon_0) \quad (72)$$

where

$$\begin{aligned}
\epsilon_0 &= \Theta \left( \tilde{\alpha}^2 \tau_v \|A^*\|_1^2 + \tau_v^2 \|A^*\|_1 \mathfrak{B}_{\mathcal{F}} + \alpha \tau_v \|A^*\|_1 \mathfrak{L}_{\mathcal{G}} \mathfrak{B}_{\mathcal{F}} \right)^2 \\
&= \tilde{\Theta} \left( \tilde{\alpha}^2 \tau_v \|A^*\|_1^2 + \alpha^2 (kp_{\mathcal{G}} \mathfrak{C}_{\mathfrak{s}}(\mathcal{G}, \mathfrak{B}_{\mathcal{F}} \|A^*\|_1))^2 (\mathfrak{B}_{\mathcal{F}} \|A^*\|_1)^3 \right. \\
&\quad \left. + \alpha^2 kp_{\mathcal{G}} \mathfrak{C}_{\mathfrak{s}}(\mathcal{G}, \mathfrak{B}_{\mathcal{F}} \|A^*\|_1) (\mathfrak{B}_{\mathcal{F}} \|A^*\|_1)^3 \|A^*\|_1^2 \right) \\
&= \tilde{\Theta} \left( \alpha^4 (p_{\mathcal{G}} \mathfrak{C}_{\mathfrak{s}}(\mathcal{G}, \mathfrak{B}_{\mathcal{F}} \|A^*\|_1))^4 (\|A^*\|_1 \mathfrak{B}_{\mathcal{F}})^6 \right)
\end{aligned} \tag{73}$$

In practice, for computational efficiency, we use the sampled adjacency matrix  $A^t$  in the learning network, so we should consider the discrepancy between the target function and the practical output

$$\begin{aligned}
\|\mathcal{H}_{n,A^*}(X) - \text{out}_n(X, A^{1t}, A^{2t})\|_2^2 &\leq \|\mathcal{H}_{n,A^*}(X) - \text{out}_n(X, A^*)\|_2^2 \\
&\quad + \|\text{out}_n(X, A^{1t}, A^{2t}) - \text{out}_n(X, A^*)\|_2^2
\end{aligned} \tag{74}$$

We have already considered  $\|\mathcal{H}_{n,A^*}(X) - \text{out}_n(X, A^*)\|_2^2$  and

$$\begin{aligned}
\|\text{out}_n(X, A^{1t}, A^{2t}) - \text{out}_n(X, A^*)\|_2 &\leq \|C\sigma(WXa_n^{1t}) - C\sigma(WXa_n^*)\|_2 \\
&\quad + \|C\sigma(V \text{out}_n^1(XA^{1t})a_n^{2t}) - C\sigma(V \text{out}_n^1(XA^*)a_n^*)\|_2
\end{aligned} \tag{75}$$

Using (27), we have

$$\|C\sigma(WXa_n^{1t}) - C\sigma(WXa_n^*)\|_2 \leq \tau_w \|Xa_n^{1t} - Xa_n^*\|_2 \leq \|\text{Err}_t\|_2 \tag{76}$$

For the above equation to hold, it requires

$$\|A^{1t} - A^*\|_1 \leq \left\| \frac{\text{Err}_t}{\tau_w} \right\|_2 \tag{77}$$

Using  $\|A^*\|_1 \leq O(1)$  and (31), we have

$$\begin{aligned}
\|C\sigma(V \text{out}_n^1(XA^{1t})a_n^{2t}) - C\sigma(V \text{out}_n^1(XA^*)a_n^*)\|_2 &\leq \tau_v \|\text{out}_n^1(XA^{1t})a_n^{2t} - \text{out}_n^1(XA^*)a_n^*\|_2 \\
&\leq \tau_v \tau_w \|A^{2t} - A^*\|_1 \leq \|\text{Err}_t\|_2
\end{aligned} \tag{78}$$

For the above equation to hold, it requires  $\|A^{2t} - A^*\|_1 \leq \left\| \frac{\text{Err}_t}{\tau_v \tau_w} \right\|_2$ .

Under assumptions of Lemma C.7, with high probability, we can ensure  $\|A^{1t} - A^*\|_2 \leq \left\| \frac{\text{Err}_t}{\tau_w} \right\|_2$ ,  $\|A^{2t} - A^*\|_1 \leq \left\| \frac{\text{Err}_t}{\tau_v \tau_w} \right\|_2$ .

Using  $\|\mathbf{W}^*\|_F \leq \tau_w/10$ ,  $\|\mathbf{V}^*\|_F \leq \tau_v/10$  and  $\epsilon \geq \text{OPT} + \epsilon_0$ , we have

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}_{n \in \mathcal{V}, (X, y_n) \sim \mathcal{D}} \|\mathcal{H}_{n,A^*}(X) - \text{out}_n(X, A^{1t}, A^{2t})\|_2^2 \leq O(\epsilon) \tag{79}$$

as long as  $T \geq \Omega\left(\frac{\tau_w^2/\eta_w + \tau_v^2/\eta_v}{\epsilon}\right)$ .

Finally, we should check  $\|\mathbf{W}_t\|_F \leq \tau_w$  and  $\|\mathbf{V}_t\|_F \leq \tau_v$  hold.

$$\frac{\|\mathbf{W}_{T_0} - \mathbf{W}^*\|_F^2}{2\eta_w T_0} + \frac{\|\mathbf{V}_{T_0} - \mathbf{V}^*\|_F^2}{2\eta_v T_0} \leq \frac{\|\mathbf{W}^*\|_F^2}{2\eta_w T_0} + \frac{\|\mathbf{V}^*\|_F^2}{2\eta_v T_0} + O(\epsilon) + \tilde{O}\left(\frac{\tau_w \|A^*\|_1}{\sqrt{T_0}}\right) \tag{80}$$

Using the relationship  $\frac{\tau_w}{\eta_w} = \frac{\tau_v}{\eta_v}$ , we have

$$\frac{\|\mathbf{W}_{T_0}\|_F^2}{\tau_w^2} + \frac{\|\mathbf{V}_{T_0}\|_F^2}{\tau_v^2} \leq \frac{4\|\mathbf{W}^*\|_F^2}{\tau_w^2} + \frac{4\|\mathbf{V}^*\|_F^2}{\tau_v^2} + 0.1 + \tilde{O}\left(\frac{\eta_w \|A^*\|_1 \sqrt{T_0}}{\tau_w}\right) \tag{81}$$

Therefore, choosing  $T = \tilde{\Theta}\left(\frac{\tau_w^2}{\|A^*\|_1 \min\{1, \epsilon^2\}}\right)$  and  $\eta_w = \tilde{\Theta}(\min\{1, \epsilon\}) \leq 0.1$ , we can ensure  $\frac{\|\mathbf{W}_{T_0}\|_F^2}{\tau_w^2} + \frac{\|\mathbf{V}_{T_0}\|_F^2}{\tau_v^2} \leq 1$ .

### C.5 Graph sampling

**Lemma C.7.** *Given a graph  $G$  with the minimum degree  $\delta(G) \geq \Omega\left(\left\|\frac{\tau_w}{Err_t}\right\|_2\right)$ , in iteration  $t$ , for the first layer  $A^{1t}$  is generated from the sampling strategy with the sampling probability  $p_{ij}^1 \leq O\left(\frac{\sqrt{d_i d_j} \|Err_t\|_2}{n_{ij} \tau_w}\right)$  and for the second layer  $A^{2t}$  is generated from the sampling strategy with the sampling probability  $p_{ij}^2 \leq O\left(\frac{\sqrt{d_i d_j} \|Err_t\|_2}{n_{ij} \tau_w \tau_v}\right)$ , we have*

$$\Pr\left[\|A^{1t} - A^*\|_1 \leq O\left(\left\|\frac{Err_t}{\tau_w}\right\|_2\right)\right] \leq e^{-\Omega(\|Err_t\|_2 \sqrt{d_i d_j} / \tau_w)} \quad (82)$$

$$\Pr\left[\|A^{2t} - A^*\|_1 \leq O\left(\left\|\frac{Err_t}{\tau_w \tau_v}\right\|_2\right)\right] \leq e^{-\Omega(\|Err_t\|_2 \sqrt{d_i d_j} / \tau_w \tau_v)} \quad (83)$$

Proof: The difference between  $A_{B_{ij}}^t$  and  $A_{B_{ij}}^*$  is

$$\Delta_{B_{ij}} = \left\|A_{B_{ij}}^t - A_{B_{ij}}^*\right\| = \sum_{i=1}^{n_{ij}} (a_{ij}^t - a_{ij}^*) \quad (84)$$

where  $n_{ij}$  is the number of elements in  $A_{B_{ij}}^*$  and  $(a_{ij}^t - a_{ij}^*)$  are iid, with  $\mu_{ij} = \mathbb{E}[\Delta_{B_{ij}}] = n_{ij} p_{ij} \frac{1}{\sqrt{d_i d_j}}$ . The Moment-generating function of  $(a_{ij}^t - a_{ij}^*)$  is

$$\begin{aligned} M_{(a_{ij}^t - a_{ij}^*)}(s) &= \mathbb{E}\left[e^{s(a_{ij}^t - a_{ij}^*)}\right] \\ &= e^{s \frac{1}{\sqrt{d_i d_j}} p_{ij}} + e^{s \cdot 0} (1 - p_{ij}) \\ &= 1 + p_{ij} \left(e^{\frac{s}{\sqrt{d_i d_j}}} - 1\right) \\ &\leq \exp\left(e^{\frac{s}{\sqrt{d_i d_j}}} - 1\right) \end{aligned} \quad (85)$$

Thus, for any  $t > 0$ , using Markov's inequality and the definition of MGF, we have

$$\begin{aligned} \mathbb{P}(\Delta_{B_{ij}} \geq k) &\leq \min_{s>0} \frac{\prod_{i=1}^{n_{ij}} M_{(a_{ij}^t - a_{ij}^*)}(s)}{e^{tk}} \\ &= \min_{t>0} \frac{e^{\mu \sqrt{d_i d_j} \left(e^{\frac{s}{\sqrt{d_i d_j}}} - 1\right)}}{e^{tk}} \end{aligned} \quad (86)$$

If  $0 \leq \delta_{ij} \leq 1$ , we plug in  $k_{ij} = (1 + \delta_{ij})\mu_{ij}$  and the optimal value of  $s_{ij} = \sqrt{d_i d_j} \ln(1 + \epsilon_{ij})$  to the above equation:

$$\mathbb{P}(\Delta_{B_{ij}} \geq (1 + \delta_{ij})\mu_{ij}) \leq \left(\frac{e^{\epsilon_{ij}}}{(1 + \epsilon_{ij})^{(1 + \epsilon_{ij})}}\right)^{\mu_{ij} \sqrt{d_i d_j}} \leq \exp\left(\frac{-\epsilon_{ij}^2 \mu_{ij} \sqrt{d_i d_j}}{3}\right) \quad (87)$$

$$\begin{aligned} (1 + \delta_{ij})^{(1 + \delta_{ij})} &= \exp[(1 + \delta_{ij}) \ln(1 + \delta_{ij})] \\ &= \exp\left(\delta_{ij} + \frac{\delta_{ij}^2}{2} - \frac{\delta_{ij}^3}{6} + o(\delta_{ij}^4)\right) \geq \exp\left(\delta_{ij} + \frac{\delta_{ij}^2}{2} - \frac{\delta_{ij}^3}{6}\right) \end{aligned} \quad (88)$$

Let  $\delta_{ij} = 1$ ,  $\mu_{ij} = \Theta(\text{Err}_t)$ , and  $d_i \geq \Omega(\frac{1}{\text{Err}_t})$ , we have  $p_{ij} \leq O(\frac{\sqrt{d_i d_j \text{Err}_t}}{n_{ij}})$ .

## C.6 Sample Complexity

**Lemma C.8.** *Given a graph  $G$  with  $|V(G)| = N$ , if the maximum degree  $\Delta(G) \leq O((N\epsilon^2)^{\frac{1}{4}})$  and sample complexity  $\Omega \geq O(\frac{\Delta(G)^2(\tau_w \|A\|_1)^2 \log N}{\epsilon^2})$ , with probability  $1 - N^{-\tau_w \|A\|_1}$ , we have*

$$\left| \mathbb{E}_{n \in \mathcal{V}, (X, y_n) \sim \mathcal{Z}} \|y_n - \text{out}_n(X, A^{1t}, A^{2t})\|_2 - \mathbb{E}_{n \in \mathcal{V}, (X, y_n) \sim \mathcal{D}} \|y_n - \text{out}_n(X, A^{1t}, A^{2t})\|_2 \right| \leq \epsilon.$$

Proof: For the set of samples  $\mathcal{Z}$  define

$$\mathbb{E}_{n \in \mathcal{V}, (X, y_n) \sim \mathcal{Z}} \|y_n - \text{out}_n(X, A^{1t}, A^{2t})\|_2 = \frac{1}{\Omega} \sum_{n=1}^{\Omega} \|y_n - \text{out}_n(X, A^{1t}, A^{2t})\|_2 \quad (89)$$

Denote the generalization error as

$$\begin{aligned} & \left| \mathbb{E}_{n \in \mathcal{V}, (X, y_n) \sim \mathcal{Z}} \|y_n - \text{out}_n(X, A^{1t}, A^{2t})\|_2 - \mathbb{E}_{n \in \mathcal{V}, (X, y_n) \sim \mathcal{D}} \|y_n - \text{out}_n(X, A^{1t}, A^{2t})\|_2 \right| \\ &= \left| \mathbb{E}_{n \in \mathcal{V}, (X, y_n) \sim \mathcal{Z}} \|\text{out}_n(X, A^{1t}, A^{2t})\|_2 - \mathbb{E}_{n \in \mathcal{V}, (X, y_n) \sim \mathcal{D}} \|\text{out}_n(X, A^{1t}, A^{2t})\|_2 \right| \end{aligned}$$

By Hoeffding's inequality and  $\|\text{out}_n(X, a_n)\|_2 \leq O(\tau_w \|A\|_1)$ , we have

$$\mathbb{E} \left[ e^{s \left| \mathbb{E}_{n \in \mathcal{V}, (X, y_n) \sim \mathcal{Z}} \|\text{out}_n(X, A^{1t}, A^{2t})\|_2 - \mathbb{E}_{n \in \mathcal{V}, (X, y_n) \sim \mathcal{D}} \|\text{out}_n(X, A^{1t}, A^{2t})\|_2 \right|} \right] \leq e^{\frac{(s\tau_w \|A\|_1)^2}{8}} \quad (90)$$

Define maximum degree of  $G$  is  $\Delta(G)$ . It is easy to know that  $\|\text{out}_n(X, A^{1t}, A^{2t})\|_2$  is dependent with at most its second order neighbor, so the maximum number of nodes related with  $\|\text{out}_n(X, A^{1t}, A^{2t})\|_2$  is  $\Delta(G)^2$ . By Lemma 7 in Shuai, we have

$$\mathbb{E} e^{s \sum_{n=1}^{\Omega} \|\text{out}_n(X, A^{1t}, A^{2t})\|_2} \leq e^{\Delta(G)^2 (s\tau_w \|A\|_1)^2 \Omega / 8} \quad (91)$$

$$\begin{aligned} \mathbb{P} \left( \left| \mathbb{E}_{n \in \mathcal{V}, (X, y_n) \sim \mathcal{Z}} \|\text{out}_n(X, A^{1t}, A^{2t})\|_2 - \mathbb{E}_{n \in \mathcal{V}, (X, y_n) \sim \mathcal{D}} \|\text{out}_n(X, A^{1t}, A^{2t})\|_2 \right| \geq \epsilon \right) \\ \leq \exp(\Delta(G)^2 (s\tau_w \|A\|_1)^2 \Omega / 8 - s\epsilon \Omega) \end{aligned} \quad (92)$$

Let  $s = \frac{4\epsilon}{\Delta(G)^2 (\tau_w \|A\|_1)^2}$  and  $\epsilon = (\tau_w \|A\|_1)^2 \sqrt{\frac{\Delta(G)^4 \log N}{\Omega}}$

$$\begin{aligned} \mathbb{P} \left( \left| \mathbb{E}_{n \in \mathcal{V}, (X, y_n) \sim \mathcal{Z}} \|\text{out}_n(X, A^{1t}, A^{2t})\|_2 - \mathbb{E}_{n \in \mathcal{V}, (X, y_n) \sim \mathcal{D}} \|\text{out}_n(X, A^{1t}, A^{2t})\|_2 \right| \geq \epsilon \right) \\ \leq \exp(-(\tau_w \|A\|_1)^2 \log N) \\ \leq N^{-\tau_w \|A\|_1} \end{aligned} \quad (93)$$

with

$$\Omega \geq O\left(\frac{\Delta(G)^2 (\tau_w \|A\|_1)^2 \log N}{\epsilon^2}\right) \quad (94)$$