
Graph Neural Networks Benefit from Structural Information Provably: A Feature Learning Perspective

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Graph neural networks (GNNs) have shown remarkable capabilities in learning
2 from graph-structured data, outperforming traditional multilayer perceptrons
3 (MLPs) in numerous graph applications. Despite these advantages, there has been
4 limited theoretical exploration into why GNNs are so effective, particularly from
5 the perspective of feature learning. This study aims to address this gap by examining
6 the role of graph convolution in feature learning theory under a specific data
7 generative model. We undertake a comparative analysis of the optimization and
8 generalization between two-layer graph convolutional networks (GCNs) and their
9 convolutional neural network (CNN) counterparts. Our findings reveal that graph
10 convolution significantly enhances the regime of low test error over CNNs. This
11 highlights a substantial discrepancy between GNNs and MLPs in terms of general-
12 ization capacity, a conclusion further supported by our empirical simulations on
13 both synthetic and real-world datasets.

14 1 Introduction

15 Graph neural networks (GNNs) have recently demonstrated remarkable capability in learning graph
16 representations, yielding superior results across various downstream tasks, such as node classifica-
17 tions [1, 2, 3], graph classifications [4, 5, 6, 7] and link predictions [8, 9, 10], etc. However, the
18 theoretical understanding of why GNNs can achieve such success is still in its infancy. Compared to
19 multilayer perceptron (MLPs), GNNs enhance representation learning with an added message passing
20 operation [11]. Take graph convolutional network (GCN) [1] as an example, it aggregates a node’s
21 attributes with those of its neighbors through a *graph convolution* operation. This operation, which
22 leverages the structural information (adjacency matrix) of graph data, forms the core distinction
23 between GNNs and MLPs. Empirical evidence from three node classification tasks, as shown in
24 Figure 1, suggests GCNs outperform MLPs. Motivated by the superior performance of GNNs, we
25 pose a critical question about graph convolution:

26 *What role does graph convolution play during gradient descent training, and what mechanism*
27 *enables a GCN to exhibit better generalization after training?*

28 Several recent studies have embarked on a theoretical exploration of graph convolution’s role in
29 GNNs. For instance, [12] considered a setting of linear classification of data generated from a
30 contextual stochastic block model [13]. Their findings indicate that graph convolution extends the
31 regime where data is linearly separable by a factor of approximately $1/\sqrt{D}$ compared to MLPs, with
32 D denoting a node’s expected degree. [14] further investigated the impact of graph convolutions in
33 multi-layer networks, showcasing improved non-linear separability. While insightful, these studies
34 assume the Bayes optimal classifier of GNNs, thereby missing a comprehensive characterization
35 of the GNNs’ optimization process. This leaves a notable gap in understanding of the optimization

36 and generalization capabilities of GNNs, a gap that existing theoretical explorations have yet to
 37 adequately address.

38 To respond to the growing demand for a comprehensive
 39 understanding of graph convolution, we delve into the feature
 40 learning analysis [15, 16]. In our study, we introduce
 41 a data generation model—termed SNM-SBM—that combines
 42 a signal-noise model [15, 17] for feature creation
 43 and a stochastic block model [18] for graph construction.
 44 Our analysis is centered on the convergence and generalization
 45 attributes of two-layer graph convolution networks
 46 (GCNs) when trained via gradient descent, compared with
 47 the established outcomes for two-layer convolutional neural
 48 networks (CNNs) as presented by [15]. While both
 49 GCNs and CNNs demonstrate to achieve near-zero training
 50 error, our study effectively sheds light on the discrepancies
 51 in their generalization abilities. We emphasize the crucial
 52 contribution of graph convolution to the enhanced performance
 of GNNs. Our study’s key contributions are as follows:

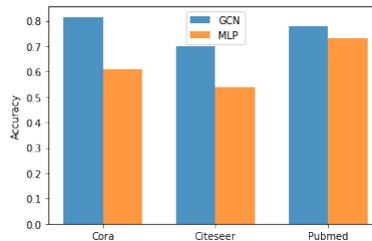


Figure 1: Performance comparison between GCN and MLP on node classification tasks.

- 53 • We establish global convergence guarantees for graph neural networks training on data drawn
 54 from SNM-SBM model by characterizing the signal learning and noise memorization in feature
 55 learning. We demonstrate that, despite the nonconvex optimization landscape, GCNs can achieve
 56 zero training error after a polynomial number of iterations.
- 57 • We further establish population loss bounds of overfitted GNN models trained by gradient descent.
 58 We show that under certain conditions on the signal-to-noise ratio, GNNs trained by gradient
 59 descent can achieve near zero test error.
- 60 • We show a contrast in the generalization of GCNs and CNNs. We identify a regime where GCNs
 61 can attain nearly zero test error, whereas the test error of CNNs is greater than a constant. This
 62 conclusion is further supported by empirical verification on synthetic and real-world datasets.

63 2 Problem Setup and Preliminary

64 **Data model** In our approach, we utilize a signal-noise model for feature generation, combined with
 65 a stochastic block model for graph structure generation. Specifically, we define the feature matrix
 66 as $\mathbf{X} \in \mathbb{R}^{n \times 2d}$, with n representing the number of samples and $2d$ being the feature dimensionality.
 67 Each feature associated with a data point is generated from a *signal-noise model* (SNM), conditional
 68 on the Rademacher random variable $y \in \{-1, 1\}$, and a latent vector $\boldsymbol{\mu} \in \mathbb{R}^d$:

$$\mathbf{x} = [\mathbf{x}^{(1)}, \mathbf{x}^{(2)}] = [y\boldsymbol{\mu}, \boldsymbol{\xi}], \quad (1)$$

69 where $\mathbf{x}^{(1)}, \mathbf{x}^{(2)} \in \mathbb{R}^d$, and $\boldsymbol{\xi} \sim \mathcal{N}(\mathbf{0}, \sigma_p^2 \cdot (\mathbf{I} - \|\boldsymbol{\mu}\|_2^{-2} \cdot \boldsymbol{\mu}\boldsymbol{\mu}^\top))$ is a Gaussian with σ_p^2 as the variance.
 70 The term $\mathbf{I} - \|\boldsymbol{\mu}\|_2^{-2} \cdot \boldsymbol{\mu}\boldsymbol{\mu}^\top$ is employed to guarantee that the noise vector is orthogonal to the
 71 signal vector $\boldsymbol{\mu}$. Moreover, we implement a stochastic block model with inter-class edge probability
 72 p and intra-class edge probability s . Specifically, the entry of adjacency matrix $\mathbf{A} = (a_{ij})_{n \times n}$
 73 is Bernoulli distributed, with $a_{ij} \sim \text{Ber}(p)$ when $y_i = y_j$, and $a_{ij} \sim \text{Ber}(s)$ when $y_i = -y_j$.
 74 The combination of a stochastic block model with the signal-noise model (1) is represented as
 75 SNM – SBM($n, p, s, \boldsymbol{\mu}, \sigma_p, d$). Note that when $p = s = 0$, SNM – SBM reduces to a SNM, and
 76 its samples are used in MLP.

77 **GCN.** Graph neural network (GNNs) fuse graph structure information and node features to learn
 78 representation of nodes. Consider a two-layer GCN f with graph convolution operation on the
 79 first layer. The output of the GCN is given by $f(\mathbf{W}, \tilde{\mathbf{x}}) = F_{+1}(\mathbf{W}_{+1}, \tilde{\mathbf{x}}) - F_{-1}(\mathbf{W}_{-1}, \tilde{\mathbf{x}})$, where
 80 $F_{+1}(\mathbf{W}_{+1}, \tilde{\mathbf{x}})$ and $F_{-1}(\mathbf{W}_{-1}, \tilde{\mathbf{x}})$ are defined as follows:

$$F_j(\mathbf{W}_j, \tilde{\mathbf{x}}) = \frac{1}{m} \sum_{r=1}^m \left[\sigma(\mathbf{w}_{j,r}^\top \tilde{\mathbf{x}}^{(1)}) + \sigma(\mathbf{w}_{j,r}^\top \tilde{\mathbf{x}}^{(2)}) \right]. \quad (2)$$

81 Here, $\tilde{\mathbf{X}} \triangleq [\tilde{\mathbf{x}}_1, \tilde{\mathbf{x}}_2, \dots, \tilde{\mathbf{x}}_n]^\top = \tilde{\mathbf{D}}^{-1} \tilde{\mathbf{A}} \mathbf{X} \in \mathbb{R}^{n \times 2d}$ with $\tilde{\mathbf{A}} = \mathbf{A} + \mathbf{I}_n$ representing the adjacency
 82 matrix with self-loop, and $\tilde{\mathbf{D}}$ is a diagonal matrix that records the degree of each node, namely,

83 $\tilde{D}_{ii} = \sum_j \tilde{A}_{ij}$. For simplicity we denote $D_i \triangleq \tilde{D}_{ii}$. Therefore, in contrast to the CNN model (6),
 84 the GCNs (2) incorporate the normalized adjacency matrix $\tilde{\mathbf{D}}^{-1}\tilde{\mathbf{A}}$, also termed as graph convolution,
 85 which serves as a pivotal component.

86 With the training data $\mathcal{S} \triangleq \{\mathbf{x}_i, y_i\}_{i=1}^n$ and \mathbf{A} drawn from SNM – SBM($n, p, s, \boldsymbol{\mu}, \sigma_p, d$), we
 87 consider to learn the network’s parameter \mathbf{W} by optimizing the cross-entropy loss function:

$$L_{\mathcal{S}}^{\text{GCN}}(\mathbf{W}) = \frac{1}{n} \sum_{i=1}^n \ell(y_i \cdot f(\mathbf{W}, \tilde{\mathbf{x}}_i)), \quad (3)$$

88 where $\ell(y \cdot f(\mathbf{W}, \mathbf{x})) = \log(1 + \exp(-f(\mathbf{W}, \mathbf{x}) \cdot y))$. The gradient descent update for the first layer
 89 weight \mathbf{W} in GCN can be expressed as:

$$\begin{aligned} \mathbf{w}_{j,r}^{(t+1)} &= \mathbf{w}_{j,r}^{(t)} - \eta \cdot \nabla_{\mathbf{w}_{j,r}} L_{\mathcal{S}}^{\text{GCN}}(\mathbf{W}^{(t)}) \\ &= \mathbf{w}_{j,r}^{(t)} - \frac{\eta}{nm} \sum_{i=1}^n \ell'_i{}^{(t)} \sigma'(\langle \mathbf{w}_{j,r}^{(t)}, \tilde{\boldsymbol{\xi}}_i \rangle) \cdot j y_i \tilde{\boldsymbol{\xi}}_i - \frac{\eta}{nm} \sum_{i=1}^n \ell'_i{}^{(t)} \sigma'(\langle \mathbf{w}_{j,r}^{(t)}, \tilde{y}_i \boldsymbol{\mu} \rangle) \cdot j \tilde{y}_i \boldsymbol{\mu}, \end{aligned} \quad (4)$$

90 where we define the loss derivative as $\ell'_i \triangleq \ell'(y_i \cdot f_i) = -\frac{\exp(-y_i \cdot f_i)}{1 + \exp(-y_i \cdot f_i)}$, “aggregated label” $\tilde{y}_i =$
 91 $D_i^{-1} \sum_{k \in \mathcal{N}(i)} y_k$ and “aggregated noise vector” $\tilde{\boldsymbol{\xi}}_i = D_i^{-1} \sum_{k \in \mathcal{N}(i)} \boldsymbol{\xi}_k$, with $\mathcal{N}(i)$ being a set
 92 that contains all the neighbor of node i . Our primary objective is to demonstrate the enhanced
 93 feature learning capabilities of GNNs in comparison to CNNs. This is achieved by examining the
 94 generalization ability of the GNN model through the lens of population loss, which can be formulated
 95 as $L_{\mathcal{D}}^{\text{GCN}}(\mathbf{W}) = \mathbb{E}_{\mathbf{x}, y \sim \mathcal{D}=\text{SNM-SBM}} \ell(y \cdot f(\mathbf{W}, \tilde{\mathbf{x}}))$.

96 In this study, our primary objective is to demonstrate the enhanced feature learning capabilities of
 97 GNNs in comparison to CNNs. This is achieved by examining the generalization ability of the
 98 GNN model through the lens of test error (population loss), which is defined based on unseen test
 99 data. Given n training data points and the corresponding graph structure, we train a GNN model.
 100 We then generate a new test data point following the SNM – SBM distribution. Its connection
 101 in the graph to the training data points are still following the stochastic block model, forming
 102 an adjacency matrix $\mathbf{A}' \in \mathbb{R}^{(n+1) \times (n+1)}$. We specifically study the population loss by taking
 103 the expectation over the randomness of the new test data, which is formulated as $L_{\mathcal{D}}^{\text{GCN}}(\mathbf{W}) =$
 104 $\mathbb{E}_{(\mathbf{x}, y, \mathbf{A}') \sim \text{SNM-SBM}} \ell(y \cdot f(\mathbf{W}, \mathbf{x}))$.

105 3 Theoretical Results

106 In this section, we introduce our key theoretical findings that explain the optimization and general-
 107 ization processes of feature learning in GCNs. Through the application of the gradient descent rule
 108 outlined in Equation (4), we observe that the gradient descent iterate $\mathbf{w}_{j,r}^{(t)}$ is a linear combination of
 109 its random initialization $\mathbf{w}_{j,r}^{(0)}$, the signal vector $\boldsymbol{\mu}$ and the noise vectors in the training data $\boldsymbol{\xi}_i^1$ for
 110 $i \in [n]$ [15]. Consequently, for $r \in [m]$, the decomposition of weight can be expressed:

$$\mathbf{w}_{j,r}^{(t)} = \mathbf{w}_{j,r}^{(0)} + j \cdot \gamma_{j,r}^{(t)} \cdot \|\boldsymbol{\mu}\|_2^{-2} \cdot \boldsymbol{\mu} + \sum_{i=1}^n \bar{\rho}_{j,r,i}^{(t)} \cdot \|\boldsymbol{\xi}_i\|_2^{-2} \cdot \boldsymbol{\xi}_i + \sum_{i=1}^n \underline{\rho}_{j,r,i}^{(t)} \cdot \|\boldsymbol{\xi}_i\|_2^{-2} \cdot \boldsymbol{\xi}_i. \quad (5)$$

111 where $\gamma_{j,r}^{(t)}$ and $\rho_{j,r,i}^{(t)} = \{\bar{\rho}_{j,r,i}^{(t)}, \underline{\rho}_{j,r,i}^{(t)}\}$ serve as coefficients. To facilitate a fine-grained analysis
 112 for the evolution of coefficients, we introduce the notations $\bar{\rho}_{j,r,i}^{(t)} \triangleq \rho_{j,r,i}^{(t)} \mathbb{1}(\rho_{j,r,i}^{(t)} \geq 0)$, $\underline{\rho}_{j,r,i}^{(t)} \triangleq$
 113 $\rho_{j,r,i}^{(t)} \mathbb{1}(\rho_{j,r,i}^{(t)} \leq 0)$. We refer to Equation (5) as the signal-noise decomposition of $\mathbf{w}_{j,r}^{(t)}$. Our analysis
 114 is based on the following assumptions:

115 **Assumption 3.1.** *Suppose that*

116 *1. The dimension d is sufficiently large: $d = \tilde{\Omega}(m^{2 \vee [4/(q-2)]} n^{4 \vee [(2q-2)/(q-2)]})$.*

¹By referring to Equation (4), we assert that the gradient descent update moves in the direction of $\tilde{\boldsymbol{\xi}}_i$ for each $i \in [n]$. Then we can apply the definition of $\tilde{\boldsymbol{\xi}}_i = D_i^{-1} \sum_{k \in \mathcal{N}(i)} \boldsymbol{\xi}_k$.

- 117 2. The size of training sample n and width of GCNs m adhere to $n, m = \Omega(\text{polylog}(d))$.
- 118 3. The learning rate η satisfies $\eta \leq \tilde{O}(\min\{\|\boldsymbol{\mu}\|_2^{-2}, \sigma_p^{-2}d^{-1}\})$.
- 119 4. The edge probability $p, s = \Omega(\sqrt{\log(n)/n})$ and $\Xi \triangleq \frac{p-s}{p+s}$ is a positive constant.
- 120 5. The standard deviation of Gaussian initialization σ_0 is chosen such that $\sigma_0 \leq$
- 121 $\tilde{O}(m^{-2/(q-2)}n^{-[1/(q-2)]\vee 1} \cdot \min\{(\sigma_p\sqrt{d/(n(p+s))})^{-1}, \Xi^{-1}\|\boldsymbol{\mu}\|_2^{-1}\})$.

122 We introduce a critical quantity called signal-to-noise ratio (SNR), which can measure the relative

123 learning speed between signal and noise, as is calculated through $\text{SNR} = \|\boldsymbol{\mu}\|_2/(\sigma_p\sqrt{d})$. To prepare

124 for our main result, we provide an effective SNR for GNNs, defined as $\text{SNR}_G = \|\boldsymbol{\mu}\|_2/(\sigma_p\sqrt{d}) \cdot$

125 $(n(p+s))^{(q-2)/(2q)}$. Given the above assumptions and definitions of SNR, we present our main

126 result for GNN as follows:

127 **Theorem 3.2.** Let $T = \tilde{\Theta}(\eta^{-1}m\sigma_0^{-(q-2)}\Xi^{-q}\|\boldsymbol{\mu}\|_2^{-q} + \eta^{-1}\epsilon^{-1}m^3\|\boldsymbol{\mu}\|_2^{-2})$. Under Assumption 3.1,

128 if $n \cdot \text{SNR}_G^q = \tilde{\Omega}(1)$, then with probability at least $1 - d^{-1}$, there exists a $0 \leq t \leq T$ such that:

- 129 • The GCN learns the signal: $\max_r \gamma_{j,r}^{(t)} = \tilde{\Omega}(1)$ for $j \in \{\pm 1\}$.
- 130 • The GCN does not memorize the noises in the training data: $\max_{j,r,i} |\rho_{j,r,i}^{(T)}| =$
- 131 $\tilde{O}(\sigma_0\sigma_p\sqrt{d/n(p+s)})$.
- 132 • The training loss converges to ϵ , i.e., $L_S^{\text{GCN}}(\mathbf{W}^{(t)}) \leq \epsilon$.
- 133 • The trained GCN achieves a small test loss: $L_D^{\text{GCN}}(\mathbf{W}^{(t)}) \leq c_1\epsilon + \exp(-c_2n^2)$.

134 where c_1 and c_2 are positive constants.

135 Theorem 3.2 reveals that, provided $n \cdot \text{SNR}_G^q = \tilde{\Omega}(1)$, the

136 GCN can learn the signal by achieving $\max_r \gamma_{j,r}^{(t)} = \tilde{\Omega}(1)$,

137 and on the other hand, the noise memorization during gra-

138 dient descent training is suppressed by $\max_{j,r,i} |\rho_{j,r,i}^{(T)}| =$

139 $\tilde{O}(\sigma_0\sigma_p\sqrt{d/n(p+s)})$, given that $\sigma_0\sigma_p\sqrt{d/n(p+s)} \ll 1$

140 according to assumption 3.1. Because the signal learned by

141 the network is large enough and much stronger than the noise

142 memory, it generalizes well to test sample. Consequently,

143 the learned neural network can achieve both small training

144 and test losses. It's worth noting that when the graph's de-

145 gree is reduced to 1, the effective SNR for GNNs converges

146 to the vanilla SNR, namely $\text{SNR}_G = \text{SNR}$. This reduces to

147 CNN, whose feature learning is established by [15].

148 Our result show that whether a graph neural network learns

149 the signal or noise depends on the SNR, and the number of

150 samples n , at the same time. According to [15], who give the

151 characterization of feature learning of CNNs, CNNs can focus on the signal learning and generalize

152 well on the unseen data when $n \cdot \text{SNR}^q = \tilde{\Omega}(1)$. On the other hand, when $n \cdot \text{SNR}^q = \tilde{O}(1)$, CNNs

153 mainly memorize the noise from data, thus achieve a large test error. To highlight the differences in

154 generalization between GNNs and CNNs, we show that, if $n \cdot \text{SNR}_G^q = \tilde{\Omega}(1)$ and $n \cdot \text{SNR}^q = \tilde{O}(1)$,

155 then the trained **GNNs achieve small test error**, given by $L_D^{\text{GCN}}(\mathbf{W}^{(t)}) = o(1)$. In contrast, the

156 trained **CNNs achieve large test error**, $L_D^{\text{CNN}}(\mathbf{W}^{(t)}) \geq C$. The first condition $n \cdot \text{SNR}_G^q = \tilde{\Omega}(1)$ is

157 by Theorem 3.2, while second the condition $n \cdot \text{SNR}^q = \tilde{\Omega}(1)$ is based on the findings of [15] for

158 CNN. As a conclusion, we clearly provide a condition that GNNs can generalize better than CNNs.

159 This observation is further visualized in Figure 2. Through the precise characterization of feature

160 learning from optimization to generalization for GNN, we have successfully demonstrated that the

161 graph neural network can gain superiority with the help of graph convolution.

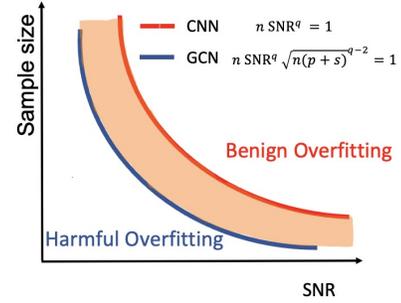


Figure 2: Illustration of performance comparison between GNN and CNN. The orange band highlights where GNN can outperform CNN.

References

- 162
- 163 [1] Thomas N Kipf and Max Welling. Semi-supervised classification with graph convolutional
164 networks. *arXiv preprint arXiv:1609.02907*, 2016.
- 165 [2] Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Lio, and Yoshua
166 Bengio. Graph attention networks. *arXiv preprint arXiv:1710.10903*, 2017.
- 167 [3] Will Hamilton, Zhitao Ying, and Jure Leskovec. Inductive representation learning on large
168 graphs. In *Advances in neural information processing systems*, pages 1024–1034, 2017.
- 169 [4] Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How powerful are graph neural
170 networks? *arXiv preprint arXiv:1810.00826*, 2018.
- 171 [5] Justin Gilmer, Samuel S Schoenholz, Patrick F Riley, Oriol Vinyals, and George E Dahl. Neural
172 message passing for quantum chemistry. *arXiv preprint arXiv:1704.01212*, 2017.
- 173 [6] Junhyun Lee, Inyeop Lee, and Jaewoo Kang. Self-attention graph pooling. *arXiv preprint
174 arXiv:1904.08082*, 2019.
- 175 [7] Hao Yuan and S. Ji. Structpool: Structured graph pooling via conditional random fields. In
176 *ICLR*, 2020.
- 177 [8] Thomas N Kipf and Max Welling. Variational graph auto-encoders. *arXiv preprint
178 arXiv:1611.07308*, 2016.
- 179 [9] Muhan Zhang and Yixin Chen. Link prediction based on graph neural networks. *Advances in
180 neural information processing systems*, 31, 2018.
- 181 [10] Ajay Kumar, Shashank Sheshar Singh, Kuldeep Singh, and Bhaskar Biswas. Link prediction
182 techniques, applications, and performance: A survey. *Physica A: Statistical Mechanics and its
183 Applications*, 553:124289, 2020.
- 184 [11] Jie Zhou, Ganqu Cui, Shengding Hu, Zhengyan Zhang, Cheng Yang, Zhiyuan Liu, Lifeng
185 Wang, Changcheng Li, and Maosong Sun. Graph neural networks: A review of methods and
186 applications. *AI open*, 1:57–81, 2020.
- 187 [12] Aseem Baranwal, Kimon Fountoulakis, and Aukosh Jagannath. Graph convolution for semi-
188 supervised classification: Improved linear separability and out-of-distribution generalization.
189 *arXiv preprint arXiv:2102.06966*, 2021.
- 190 [13] Yash Deshpande, Andrea Montanari, Elchanan Mossel, and Subhabrata Sen. Contextual
191 stochastic block models. *arXiv preprint arXiv:1807.09596*, 2018.
- 192 [14] Aseem Baranwal, Kimon Fountoulakis, and Aukosh Jagannath. Effects of graph convolutions in
193 multi-layer networks. In *The Eleventh International Conference on Learning Representations*,
194 2023.
- 195 [15] Yuan Cao, Zixiang Chen, Mikhail Belkin, and Quanquan Gu. Benign overfitting in two-layer
196 convolutional neural networks. *arXiv preprint arXiv:2202.06526*, 2022.
- 197 [16] Zeyuan Allen-Zhu and Yuanzhi Li. Feature purification: How adversarial training performs
198 robust deep learning. In *2021 IEEE 62nd Annual Symposium on Foundations of Computer
199 Science (FOCS)*, pages 977–988. IEEE, 2022.
- 200 [17] Zeyuan Allen-Zhu and Yuanzhi Li. Towards understanding ensemble, knowledge distillation
201 and self-distillation in deep learning. *arXiv preprint arXiv:2012.09816*, 2020.
- 202 [18] Emmanuel Abbe, Afonso S Bandeira, and Georgina Hall. Exact recovery in the stochastic block
203 model. *IEEE Transactions on information theory*, 62(1):471–487, 2015.
- 204 [19] Difan Zou, Yuan Cao, Yuanzhi Li, and Quanquan Gu. Understanding the generalization of
205 adam in learning neural networks with proper regularization. *arXiv preprint arXiv:2108.11371*,
206 2021.

- 207 [20] Ruoqi Shen, Sebastien Bubeck, and Suriya Gunasekar. Data augmentation as feature ma-
208 nipulation. In *International Conference on Machine Learning*, pages 19773–19808. PMLR,
209 2022.
- 210 [21] Niladri S Chatterji and Philip M Long. Finite-sample analysis of interpolating linear classifiers
211 in the overparameterized regime. *J. Mach. Learn. Res.*, 22:129–1, 2021.
- 212 [22] Zhengdao Chen, Xiang Li, and Joan Bruna. Supervised community detection with line graph
213 neural networks. *arXiv preprint arXiv:1705.08415*, 2017.
- 214 [23] Yao Ma, Xiaorui Liu, Neil Shah, and Jiliang Tang. Is homophily a necessity for graph neural
215 networks? *arXiv preprint arXiv:2106.06134*, 2021.
- 216 [24] Si Zhang, Hanghang Tong, Jiejun Xu, and Ross Maciejewski. Graph convolutional networks: a
217 comprehensive review. *Computational Social Networks*, 6(1):1–23, 2019.
- 218 [25] Xiangnan He, Kuan Deng, Xiang Wang, Yan Li, Yongdong Zhang, and Meng Wang. Lightgcn:
219 Simplifying and powering graph convolution network for recommendation. In *Proceedings of
220 the 43rd International ACM SIGIR conference on research and development in Information
221 Retrieval*, pages 639–648, 2020.
- 222 [26] Felix Wu, Amauri Souza, Tianyi Zhang, Christopher Fifty, Tao Yu, and Kilian Weinberger.
223 Simplifying graph convolutional networks. In *International conference on machine learning*,
224 pages 6861–6871. PMLR, 2019.
- 225 [27] Kun Wang, Guohao Li, Shilong Wang, Guibin Zhang, Kai Wang, Yang You, Xiaojiang Peng,
226 Yuxuan Liang, and Yang Wang. The snowflake hypothesis: Training deep gnn with one node
227 one receptive field. *arXiv preprint arXiv:2308.10051*, 2023.
- 228 [28] Keyulu Xu, Mozhi Zhang, Jingling Li, Simon S Du, Ken-ichi Kawarabayashi, and Stefanie
229 Jegelka. How neural networks extrapolate: From feedforward to graph neural networks. *arXiv
230 preprint arXiv:2009.11848*, 2020.
- 231 [29] Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and
232 generalization in neural networks. *Advances in neural information processing systems*, 31,
233 2018.
- 234 [30] Wei Huang, Yayong Li, Weitao Du, Richard Yi Da Xu, Jie Yin, Ling Chen, and Miao Zhang.
235 Towards deepening graph neural networks: A gntk-based optimization perspective. *arXiv
236 preprint arXiv:2103.03113*, 2021.
- 237 [31] Qimai Li, Zhichao Han, and Xiao-Ming Wu. Deeper insights into graph convolutional networks
238 for semi-supervised learning. In *Thirty-Second AAAI conference on artificial intelligence*, 2018.
- 239 [32] Yifan Hou, Jian Zhang, James Cheng, Kaili Ma, Richard TB Ma, Hongzhi Chen, and Ming-
240 Chang Yang. Measuring and improving the use of graph information in graph neural networks.
241 *arXiv preprint arXiv:2206.13170*, 2022.
- 242 [33] Chenxiao Yang, Qitian Wu, Jiahua Wang, and Junchi Yan. Graph neural networks are inherently
243 good generalizers: Insights by bridging gnns and mlps. *arXiv preprint arXiv:2212.09034*, 2022.
- 244 [34] Simon S Du, Kangcheng Hou, Russ R Salakhutdinov, Barnabas Poczos, Ruosong Wang, and
245 Keyulu Xu. Graph neural tangent kernel: Fusing graph neural networks with graph kernels.
246 *Advances in neural information processing systems*, 32, 2019.
- 247 [35] Mahalakshmi Sabanayagam, Pascal Esser, and Debarghya Ghoshdastidar. Representation power
248 of graph convolutions: Neural tangent kernel analysis. *arXiv preprint arXiv:2210.09809*, 2022.
- 249 [36] Jimmy Ba, Murat A Erdogdu, Taiji Suzuki, Zhichao Wang, Denny Wu, and Greg Yang. High-
250 dimensional asymptotics of feature learning: How one gradient step improves the representation.
251 *arXiv preprint arXiv:2205.01445*, 2022.
- 252 [37] Greg Yang and Edward J Hu. Feature learning in infinite-width neural networks. *arXiv preprint
253 arXiv:2011.14522*, 2020.

- 254 [38] Zixin Wen and Yuanzhi Li. Toward understanding the feature learning process of self-supervised
 255 contrastive learning. In *International Conference on Machine Learning*, pages 11112–11122.
 256 PMLR, 2021.
- 257 [39] Alexandru Damian, Jason Lee, and Mahdi Soltanolkotabi. Neural networks can learn represen-
 258 tations with gradient descent. In *Conference on Learning Theory*, pages 5413–5452. PMLR,
 259 2022.
- 260 [40] Difan Zou, Yuan Cao, Yuanzhi Li, and Quanquan Gu. The benefits of mixup for feature learning.
 261 *arXiv preprint arXiv:2303.08433*, 2023.
- 262 [41] Yongqiang Chen, Wei Huang, Kaiwen Zhou, Yatao Bian, Bo Han, and James Cheng. To-
 263 wards understanding feature learning in out-of-distribution generalization. *arXiv preprint*
 264 *arXiv:2304.11327*, 2023.
- 265 [42] Xuran Meng, Yuan Cao, and Difan Zou. Per-example gradient regularization improves learning
 266 signals from noisy data. *arXiv preprint arXiv:2303.17940*, 2023.
- 267 [43] Samy Jelassi, Michael Sander, and Yuanzhi Li. Vision transformers provably learn spatial
 268 structure. *Advances in Neural Information Processing Systems*, 35:37822–37836, 2022.
- 269 [44] Yiwen Kou, Zixiang Chen, Yuanzhou Chen, and Quanquan Gu. Benign overfitting for two-layer
 270 relu networks. *arXiv preprint arXiv:2303.04145*, 2023.

271 A Appendix for Problem Setup

272 **Notations** We use lower bold-faced letters for vectors, upper bold-faced letters for matrices, and
 273 non-bold-faced letters for scalars. For a vector \mathbf{v} , its ℓ_2 -norm is denoted as $\|\mathbf{v}\|_2$. For a matrix \mathbf{A} ,
 274 we use $\|\mathbf{A}\|_2$ to denote its spectral norm and $\|\mathbf{A}\|_F$ for its Frobenius norm. We employ standard
 275 asymptotic notations such as $O(\cdot)$, $o(\cdot)$, $\Omega(\cdot)$, and $\Theta(\cdot)$ to describe the limiting behavior. We use
 276 $\tilde{O}(\cdot)$, $\tilde{\Omega}(\cdot)$, and $\tilde{\Theta}(\cdot)$ to hide logarithmic factors in these notations respectively. Moreover, we
 277 denote $a_n = \text{poly}(b_n)$ if $a_n = O((b_n)^p)$ for some positive constant p and $a_n = \text{polylog}(b_n)$ if
 278 $a_n = \text{poly}(\log(b_n))$. Lastly, sequences of integers are denoted as $[m] = \{1, 2, \dots, m\}$.

279 The signal-noise model we have adopted is inspired by the structure of an image composed of multiple
 280 patches, where we consider a two-patch model for simplicity. The first patch $\mathbf{x}^{(1)}$, represented by
 281 the signal vector, corresponds to the target in an image. The second patch $\mathbf{x}^{(2)}$, represented by the
 282 noise vector, corresponds to the background. It’s worth mentioning that a series of recent works
 283 [17, 15, 19, 20] have explored similar signal-noise models to illustrate the feature learning process of
 284 neural networks.

285 **CNN.** We introduce a two-layer CNN model, denoted as f , which utilizes a non-linear activation
 286 function, $\sigma(\cdot)$. Specifically, we employ a polynomial ReLU activation function defined as $\sigma(z) =$
 287 $\max\{0, z\}^q$, where $q > 2$ is a hyperparameter. Mathematically, given the input data \mathbf{x} , the CNN’s
 288 output is represented as $f(\mathbf{W}, \mathbf{x}) = F_{+1}(\mathbf{W}_{+1}, \mathbf{x}) - F_{-1}(\mathbf{W}_{-1}, \mathbf{x})$, where $F_{+1}(\mathbf{W}_{+1}, \mathbf{x})$ and
 289 $F_{-1}(\mathbf{W}_{-1}, \mathbf{x})$ are defined as follows:

$$F_j(\mathbf{W}_j, \mathbf{x}) = \frac{1}{m} \sum_{r=1}^m \left[\sigma(\mathbf{w}_{j,r}^\top \mathbf{x}^{(1)}) + \sigma(\mathbf{w}_{j,r}^\top \mathbf{x}^{(2)}) \right], \quad (6)$$

290 where m is the width of hidden layer, the second layer parameters are fixed as either $+1$ or -1 ,
 291 and $\mathbf{w}_{j,r} \in \mathbb{R}^d$ refers to the weight of the first layer’s r -th. The symbol \mathbf{W} collectively represents
 292 the model’s weights. Moreover, each weight in the first layer is initialized from a random draw
 293 of a Gaussian random variable, $\mathbf{w}_{j,r} \sim \mathcal{N}(\mathbf{0}, \sigma_0^2 \cdot \mathbf{I}_{d \times d})$ for all $r \in [m]$ and $j \in \{-1, 1\}$, with σ_0
 294 regulating the initialization magnitude for the first layer’s weight.

295 Upon receiving training data $\mathcal{S} \triangleq \{\mathbf{x}_i, y_i\}_{i=1}^n$ drawn from SNM – SBM($n, p = 0, s = 0, \boldsymbol{\mu}, \sigma_p, d$),
 296 we aim to learn the parameter \mathbf{W} by minimizing the empirical cross-entropy loss function:

$$L_{\mathcal{S}}^{\text{CNN}}(\mathbf{W}) = \frac{1}{n} \sum_{i=1}^n \ell(y_i \cdot f(\mathbf{W}, \mathbf{x}_i)), \quad (7)$$

297 where $\ell(y \cdot f(\mathbf{W}, \mathbf{x})) = \log(1 + \exp(-f(\mathbf{W}, \mathbf{x}) \cdot y))$. The update rule for the gradient descent used
 298 in the CNN is then given as:

$$\begin{aligned} \mathbf{w}_{j,r}^{(t+1)} &= \mathbf{w}_{j,r}^{(t)} - \eta \cdot \nabla_{\mathbf{w}_{j,r}} L_S^{\text{CNN}}(\mathbf{W}^{(t)}) \\ &= \mathbf{w}_{j,r}^{(t)} - \frac{\eta}{nm} \sum_{i=1}^n \ell_i^{(t)} \cdot \sigma'(\langle \mathbf{w}_{j,r}^{(t)}, \boldsymbol{\xi}_i \rangle) \cdot jy_i \boldsymbol{\xi}_i - \frac{\eta}{nm} \sum_{i=1}^n \ell_i^{(t)} \cdot \sigma'(\langle \mathbf{w}_{j,r}^{(t)}, y_i \boldsymbol{\mu} \rangle) \cdot j\boldsymbol{\mu}, \end{aligned} \quad (8)$$

299 where we define the loss derivative as $\ell_i' \triangleq \ell'(y_i \cdot f_i) = -\frac{\exp(-y_i \cdot f_i)}{1 + \exp(-y_i \cdot f_i)}$. It’s important to clarify that
 300 the model we use for the MLP part is a CNN. We categorize it as an MLP for comparison purposes.

301 B Remark on Assumption 3.1

302 **Remark B.1.** (1) The requirement for the dimension d ensures that the learning process operates in
 303 a suitably over-parameterized environment [21, 15] when the second layer remains fixed. (2) It’s
 304 necessary for the sample size and neural network width to be at least polylogarithmic in the dimension
 305 d . This condition ensures certain statistical properties of the training data and weight initialization
 306 hold with a probability of at least $1 - d^{-1}$. (3) The condition on η is to ensure that gradient descent
 307 can effectively minimize the training loss. (4) The assumption regarding edge probability guarantees
 308 a sufficient level of concentration in the degree and an adequate display of homophily of graph data.
 309 (5) Lastly, the conditions imposed on initialization strength σ_0 are intended to guarantee that the
 310 training loss can effectively converge to a sufficiently small value and to discern the differential
 311 learning speed between signal and noise.

312 C Related Work

313 **Role of Graph Convolution in GNNs.** Enormous empirical studies of various GNNs models with
 314 graph convolution [22, 23, 24, 25, 26, 27] have been demonstrating that graph convolutions can
 315 enhance the performance of traditional classification methods, such as a multi-layer perceptron (MLP).
 316 Towards theoretically understanding the role of graph convolution, [28] identify conditions under
 317 which MLPs and GNNs extrapolate, thereby highlighting the superiority of GNNs for extrapolation
 318 problems. Their theoretical analysis leveraged the concept of the over-parameterized networks and
 319 the neural tangent kernel [29]. [30] use a similar approach to examine the role of graph convolution
 320 in deep GNNs within a node classification setting. They discover that excessive graph convolution
 321 layers can hamper the optimization and generalization of GNNs, corroborating the well-known
 322 over-smoothing issue in deep GNNs [31]. Another work by [32] propose two smoothness metrics
 323 to measure the quantity and quality of information derived from graph data, along with a novel
 324 attention-based framework. Some recent works [12, 14, 23] have demonstrated that graph convolution
 325 broadens the regime in which a multi-layer network can classify nodes, compared to methods that
 326 do not utilize the graph structure, especially when the graph is dense and exhibits homophily. [33]
 327 attribute the major performance gains of GNNs to their inherent generalization capability through
 328 graph neural tangent kernel (GNTK) and extrapolation analysis. As for neural network theory, these
 329 works either gleaned insights from GNTK [34, 30, 35] or studied the role of graph convolution within
 330 a linear neural network setting. Unlike them, our work is beyond NTK and investigates a more
 331 realistic setting concerning the convergence and generalization of neural networks in terms of feature
 332 learning.

333 **Feature Learning in Neural Networks.** This work builds upon a growing body of research on how
 334 neural networks learn features. [17] formulated a theory illustrating that when data possess a “multi-
 335 view” feature, ensembles of independently trained neural networks can demonstrably improve test
 336 accuracy. Further, [16] demonstrated that adversarial training can purge certain small dense mixtures
 337 from the hidden weights during the training process of a neural network, thus refining the hidden
 338 weights. [36] established that the initial gradient update contains a rank-1 ‘spike’, which leads to an
 339 alignment between the first-layer weights and the linear component feature of the teacher model. [15]
 340 investigated the benign overfitting phenomenon in training a two-layer convolutional neural network
 341 (CNN), illustrating that under certain conditions related to the signal-to-noise ratio, a two-layer CNN
 342 trained by gradient descent can achieve exceedingly low test loss through feature learning. Alongside
 343 related works [37, 19, 38, 39, 40, 41, 42, 43, 44], all these studies have highlighted the existence of

344 feature learning in neural networks during gradient descent training, forming a critical line of inquiry
 345 that this work continues to explore.

346 D Conclusion and Limitations

347 This paper utilizes a signal-noise decomposition to study the signal learning and noise memorization
 348 process in training a two-layer GCN. We provide specific conditions under which a GNN will
 349 primarily concentrate on signal learning, thereby achieving low training and testing errors. Our
 350 results theoretically demonstrate that GCNs, by leveraging structural information, outperform CNNs
 351 in terms of generalization ability across a broader benign regime. As a pioneering work that studies
 352 feature learning of GNNs, our theoretical framework is constrained to examining the role of graph
 353 convolution within a specific two-layer GCN and a certain data generalization model. In fact, the
 354 feature learning of a neural network can be influenced by a myriad of other factors, such as the depth
 355 of GNN, activation function, optimization algorithm, and data model [44, 19, 40]. Future work can
 356 extend our framework to consider the influence of a wider array of factors on feature learning within
 357 GCNs.

358 E Proof Sketches

359 In this section, we present proof sketches inspired by the study of feature learning in CNNs [15]. This
 360 foundation allows us to extend and adapt these concepts to a novel context for GNNs. We discuss
 361 the primary challenges encountered during the study of GNN, and illustrate the key techniques we
 362 employed in our proofs to overcome these challenges. These main techniques are elaborated in the
 363 following sections, and detailed proofs can be found in the appendix.

364 E.1 Iterative of coefficients under graph convolution

365 To analyze the feature learning process of graph neural networks during gradient descent training, we
 366 introduce an iterative methodology, based on the signal-noise decomposition in decomposition (5)
 367 and gradient descent update (4). The following lemma offers us a means to monitor the iteration of
 368 the signal learning and noise memorization under graph convolution:

369 **Lemma E.1.** *The coefficients $\gamma_{j,r}^{(t)}$, $\bar{\rho}_{j,r,i}^{(t)}$, $\underline{\rho}_{j,r,i}^{(t)}$ in decomposition (5) adhere to the following equa-*
 370 *tions:*

$$\gamma_{j,r}^{(0)}, \bar{\rho}_{j,r,i}^{(0)}, \underline{\rho}_{j,r,i}^{(0)} = 0, \quad (9)$$

$$\gamma_{j,r}^{(t+1)} = \gamma_{j,r}^{(t)} - \frac{\eta}{nm} \cdot \sum_{i=1}^n \ell_i^{(t)} \sigma'(\langle \mathbf{w}_{j,r}^{(t)}, \tilde{y}_i \boldsymbol{\mu}_i \rangle) y_i \tilde{y}_i \|\boldsymbol{\mu}\|_2^2, \quad (10)$$

$$\bar{\rho}_{j,r,i}^{(t+1)} = \bar{\rho}_{j,r,i}^{(t)} - \frac{\eta}{nm} \cdot \sum_{k \in \mathcal{N}(i)} D_k^{-1} \cdot \ell_k^{(t)} \cdot \sigma'(\langle \mathbf{w}_{j,r}^{(t)}, \tilde{\boldsymbol{\xi}}_k \rangle) \cdot \|\boldsymbol{\xi}_i\|_2^2 \cdot \mathbb{1}(y_k = j), \quad (11)$$

$$\underline{\rho}_{j,r,i}^{(t+1)} = \underline{\rho}_{j,r,i}^{(t)} + \frac{\eta}{nm} \cdot \sum_{k \in \mathcal{N}(i)} D_k^{-1} \cdot \ell_k^{(t)} \cdot \sigma'(\langle \mathbf{w}_{j,r}^{(t)}, \tilde{\boldsymbol{\xi}}_k \rangle) \cdot \|\boldsymbol{\xi}_i\|_2^2 \cdot \mathbb{1}(y_k = -j). \quad (12)$$

371 Lemma E.1 simplifies the analysis of the feature learning in GCNs by reducing it to the examination
 372 of the discrete dynamical system expressed by Equations (10 - 12). Our proof strategy emphasizes an
 373 in-depth evaluation of the coefficient values $\gamma_{j,r}^{(t)}$, $\bar{\rho}_{j,r,i}^{(t)}$, $\underline{\rho}_{j,r,i}^{(t)}$ throughout the training. Note that graph
 374 convolution aggregates information from neighboring nodes to the central node, which often leads to
 375 the loss of statistical stability for the aggregated noise vectors and labels. To overcome this challenge,
 376 we utilize a dense graph input, achieved by setting the edge probability as stated in 3.1.

377 E.2 A two-phase dynamics analysis

378 We then provide a two-stage dynamics analysis based on the behavior of loss derivative to track the
 379 trajectory of coefficients for signal learning and noise memorization:

380 **Stage 1.** Intuitively, the initial neural network weights are small enough so that the neural network
 381 at initialization has constant level cross-entropy loss derivatives on all the training data: $\ell_i^{(0)} =$
 382 $\ell'[y_i \cdot f(\mathbf{W}^{(0)}, \tilde{\mathbf{x}}_i)] = \Theta(1)$ for all $i \in [n]$. This is guaranteed under Condition 3.1 on σ_0 . Motivated
 383 by this, the dynamics of the coefficients in Equations (10 - 12) can be greatly simplified by replacing
 384 the $\ell_i^{(t)}$ factors by their constant upper and lower bounds. The following lemma summarizes our
 385 main conclusion at stage 1 for signal learning:

386 **Lemma E.2.** *Under the same conditions as Theorem 3.2, there exists $T_1 =$*
 387 *$\tilde{O}(\eta^{-1} m \sigma_0^{2-q} \Xi^{-q} \|\boldsymbol{\mu}\|_2^{-q})$ such that $\max_r \gamma_{j,r}^{(T_1)} = \Omega(1)$ for $j \in \{\pm 1\}$, and $|\rho_{j,r,i}^{(t)}| =$*
 388 *$O\left(\sigma_0 \sigma_p \sqrt{d} / \sqrt{n(p+s)}\right)$ for all $j \in \{\pm 1\}$, $r \in [m]$, $i \in [n]$ and $0 \leq t \leq T_1$.*

389 The proof can be found in Appendix I.1. Lemmas E.2 leverages the period of training when the
 390 derivatives of the loss function are of a constant order. It’s important to note that graph convolution
 391 plays a significant role in diverging the learning speed between signal learning and noise memorization
 392 in this first stage. Note that graph convolution can potentially cause unstable iterative dynamics of
 393 coefficients during the feature learning process. To mitigate this issue, we introduce “homophily” by
 394 setting $p > s$, which helps in stabilizing the coefficient iterations.

395 Originally, the learning speeds are roughly determined by $\|\boldsymbol{\mu}\|_2$ and $\|\boldsymbol{\xi}\|_2$ respectively without graph
 396 convolution [15]. Instead, with graph convolution, the learning speeds are approximately determined
 397 by $|\tilde{y}|\|\boldsymbol{\mu}\|_2$ and $\|\tilde{\boldsymbol{\xi}}\|_2$ respectively. Here, $|\tilde{y}|\|\boldsymbol{\mu}\|_2$ is close to $\|\boldsymbol{\mu}\|_2$, but $\|\tilde{\boldsymbol{\xi}}\|_2$ is much smaller than
 398 $\|\boldsymbol{\xi}\|_2$ (see Figure 6 for an illustration). This means that graph convolution can slow down noise
 399 memorization, thus enabling GNNs to focus more on signal learning.

400 **Stage 2.** Building on the results from the first stage, we then move to the second stage of the training
 401 process. In this stage, the loss derivatives are no longer constant, and we demonstrate that the training
 402 error can be minimized to an arbitrarily small value. Importantly, the scale differences established
 403 during the first stage of learning continue to be maintained throughout the second stage:

404 **Lemma E.3.** *Under the same conditions as Theorem 3.2, for any $t \in [T_1, T]$, it holds that*
 405 *$\max_r \gamma_{j,r}^{(T_1)} \geq 2, \forall j \in \{\pm 1\}$ and $|\rho_{j,r,i}^{(t)}| \leq \sigma_0 \sigma_p \sqrt{d} / (n(p+s))$ for all $j \in \{\pm 1\}$, $r \in [m]$*
 406 *and $i \in [n]$. Moreover, we have $L_S^{\text{GCN}}(\mathbf{W}^{(t)}) \leq \epsilon$.*

407 Lemma E.3 presents two primary outcomes. Firstly, throughout this training phase, it ensures that the
 408 coefficients of noise vectors, denoted as $\rho_{j,r,i}^{(t)}$, retain a significantly small value while coefficients of
 409 feature vector, denoted as $\gamma_{j,r}^{(t)}$ can achieve large value. Furthermore, it offers a convergence for GNN,
 410 showing the training loss will tend to receive an arbitrarily small value.

411 E.3 Test error analysis

412 Finally, it is a challenge for the generalization analysis of graph neural networks. To address this
 413 issue, we introduce an expectation over the distribution for a single data point. We consider a new
 414 data point (\mathbf{x}, y) drawn from the distribution SNM-SBM. The lemma below further gives an upper
 415 bound on the test loss of GNNs post-training:

416 **Lemma E.4.** *Let T be defined in Theorem 3.2. Under the same conditions as Theorem 3.2, for any*
 417 *$t \leq T$ with $L_S^{\text{GCN}}(\mathbf{W}^{(t)}) \leq 1$, it holds that $L_D^{\text{GCN}}(\mathbf{W}^{(t)}) \leq c_1 \cdot L_S^{\text{GCN}}(\mathbf{W}^{(t)}) + \exp(-c_2 n^2)$.*

418 The proof is presented in the appendix. Lemma E.4 demonstrates that GNNs achieve a small test
 419 error (*benign overfitting*) and completes the last step of feature learning theory.

420 F Experiments

421 In this section, we validate our theoretical findings through numerical simulations using synthetic
 422 data, specifically generated according to the SNM-SBM model. We set the signal vector, $\boldsymbol{\mu}$, to
 423 drawn from a standard normal distribution $\mathcal{N}(\mathbf{0}, \mathbf{I})$. The noise vector, $\boldsymbol{\xi}$, is drawn from a Gaussian
 424 distribution $\mathcal{N}(\mathbf{0}, \sigma_p^2 \mathbf{I})$. We train a two-layer CNN defined as per equation (6) and a two-layer GNN
 425 as per equation (2) with polynomial ReLU $q = 3$. We used the gradient descent method with a

426 learning rate of $\eta = 0.03$. The primary task we focused on was node classification, where the goal
 427 was to predict the class labels of nodes in a graph.

428 **Feature learning dynamics.** Firstly, we display the training loss, test loss, training accuracy, and
 429 test accuracy for both the CNN and GNN in Figure 3. In this case, we further set the training data size
 430 to $n = 250$, input dimension to $d = 500$, noise strength to $\sigma_p = 20$, and edge probability to $p = 0.5$,
 431 $s = 0.08$. We observe that both the GNN and CNN can achieve zero training error. However, while
 432 the GNN obtains nearly zero test error, the CNN fails to generalize effectively to the test set. This
 433 simulation result serves to validate our theoretical results in Theorem 3.2.

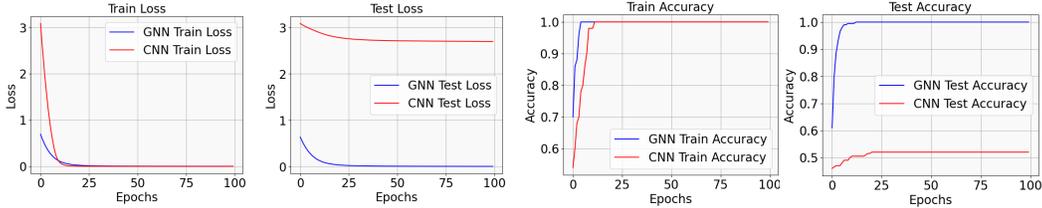


Figure 3: Training loss, testing loss, training accuracy, and testing accuracy for both CNN and GNN over a span of 100 training epochs.

434 **Verification via real-world data.** We conducted an experiment using real-world data, specifically
 435 by replacing the synthetic feature with MNIST input features. We select numbers 1 and 2 from the ten
 436 digital numbers, and applied both CNN and GNN models as described in our paper. Detailed results
 437 and visualizations can be found in the Figure 4. The results were consistent with our theoretical
 438 conclusions, reinforcing the insights derived from our analysis. We believe that this experiment adds
 439 a valuable dimension to our work, bridging the gap between theory and practice.

440 **Phase diagram.** We then explore a range of Signal-to-Noise Ratios (SNRs) from 0.045 to 0.98,
 441 and a variety of sample sizes, n , ranging from 200 to 7200. Based on our results, we train the neural
 442 network for 200 steps for each combination of SNR and sample size n . After training, we calculate
 443 the test accuracy for each run. The results are presented as a heatmap in Figure 5. Compared to
 444 CNNs, GCNs demonstrate a perfect accuracy score of 1 across a more extensive range in the SNR
 445 and n plane, indicating that GNNs have a broader *benign overfitting* regime with high test accuracy.
 446 This further validates our theoretical findings.

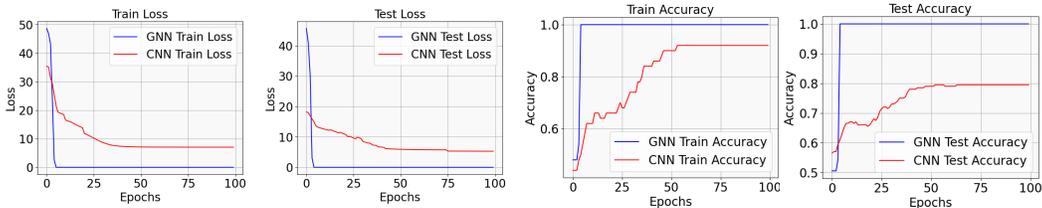


Figure 4: The verification of our theoretical result with a real-world data. The input feature is form MNIST dataset, where we select number 1 and 2 as two classes. The graph structure is sampled form stochastic block model. We show the training loss, testing loss, training accuracy, and testing accuracy for both CNN and GNN over a span of 100 training epochs. The results confirm the benefit of GNN over CNN on the real world dataset.

447 G Preliminary Lemmas

448 In this section, we present preliminary lemmas which form the foundation for the proofs to be detailed
 449 in the subsequent sections. The proof will be developed after the lemmas presented.

450 G.1 Preliminary Lemmas without Graph Convolution

451 In this section, we introduce necessary lemmas that will be used in the analysis without graph
 452 convolution, following the study of feature learning in CNN [15]. In particular, Lemma G.1 states

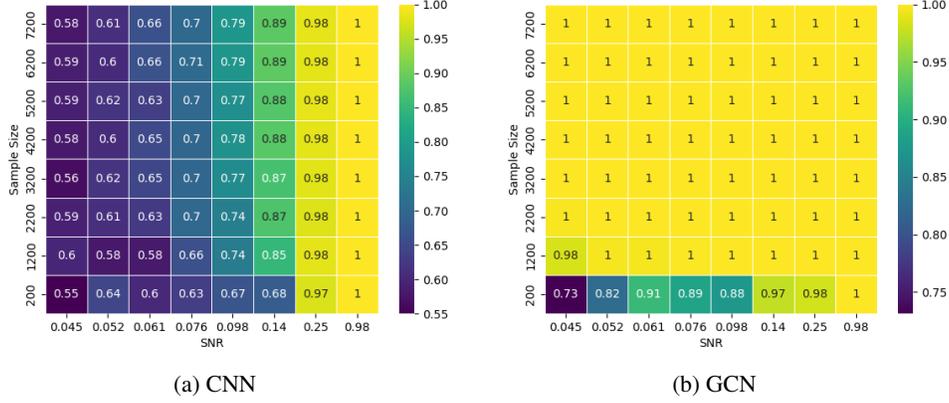


Figure 5: Test accuracy heatmap for CNNs and GCNs after training.

453 that noise vectors are “almost orthogonal” to each other and Lemma G.2 indicates that random
 454 initialization results in a controllable inner product between the weights at initialization and the data
 455 vectors.

456 **Lemma G.1.** [15] Suppose that $\delta > 0$ and $d = \Omega(\log(4n/\delta))$. Then with probability at least $1 - \delta$,

$$\begin{aligned} \sigma_p^2 d/2 &\leq \|\xi_i\|_2^2 \leq 3\sigma_p^2 d/2, \\ |\langle \xi_i, \xi_{i'} \rangle| &\leq 2\sigma_p^2 \cdot \sqrt{d \log(4n^2/\delta)}, \end{aligned}$$

457 for all $i, i' \in [n]$.

458 **Lemma G.2.** [15] Suppose that $d = \Omega(\log(nm/\delta))$, $m = \Omega(\log(1/\delta))$. Then with probability at
 459 least $1 - \delta$,

$$\begin{aligned} |\langle \mathbf{w}_{j,r}^{(0)}, \boldsymbol{\mu} \rangle| &\leq \sqrt{2 \log(8m/\delta)} \cdot \sigma_0 \|\boldsymbol{\mu}\|_2, \\ |\langle \mathbf{w}_{j,r}^{(0)}, \xi_i \rangle| &\leq 2\sqrt{\log(8mn/\delta)} \cdot \sigma_0 \sigma_p \sqrt{d}, \end{aligned}$$

460 for all $r \in [m]$, $j \in \{\pm 1\}$ and $i \in [n]$. Moreover,

$$\begin{aligned} \sigma_0 \|\boldsymbol{\mu}\|_2/2 &\leq \max_{r \in [m]} j \cdot \langle \mathbf{w}_{j,r}^{(0)}, \boldsymbol{\mu} \rangle \leq \sqrt{2 \log(8m/\delta)} \cdot \sigma_0 \|\boldsymbol{\mu}\|_2, \\ \sigma_0 \sigma_p \sqrt{d}/4 &\leq \max_{r \in [m]} j \cdot \langle \mathbf{w}_{j,r}^{(0)}, \xi_i \rangle \leq 2\sqrt{\log(8mn/\delta)} \cdot \sigma_0 \sigma_p \sqrt{d}, \end{aligned}$$

461 for all $j \in \{\pm 1\}$ and $i \in [n]$.

462 G.2 Preliminary Lemmas on Graph Properties

463 We now introduce important lemmas that are critical to our analysis. The key idea to ensure a
 464 relatively dense graph. In a sparser graph, the concentration properties of graph degree (Lemma
 465 G.3), the graph convoluted label (G.4), the graph convoluted noise vector (Lemma G.7 and Lemma
 466 G.5) are no longer guaranteed. This lack of concentration affects the behavior of coefficients during
 467 gradient descent training, leading to deviations from our current main results.

468 **Lemma G.3** (Degree concentration). Let $p, s = \Omega\left(\sqrt{\frac{\log(n/\delta)}{n}}\right)$ and $\delta > 0$, then with probability at
 469 least $1 - \delta$, we have

$$n(p + s)/4 \leq D_i \leq 3n(p + s)/4.$$

470 *Proof.* It is known that the degrees are sums of Bernoulli random variables.

$$D_i = 1 + \sum_{j \neq i}^n a_{ij},$$

471 where $a_{ij} = [\mathbf{A}]_{ij}$. Hence, by the Hoeffding's inequality, with probability at least $1 - \delta/n$

$$|D_i - \mathbb{E}[D_i]| < \sqrt{\log(n/\delta)(n-1)}.$$

472 Note that $a_{ii} = 1$ is a fixed value, which means that it is not a random variable, thus the denominator
473 in the exponential part is $n - 1$ instead of n . Now we calculate the expectation of degree:

$$\mathbb{E}[D_{ii}] = 1 + \frac{n}{2}s + \left(\frac{n}{2} - 1\right)p = n(p+s)/2 + 1 - p,$$

474 then we have

$$|D_i - n(p+s)/2 + 1 - p| \leq \sqrt{n \log(n/\delta)}.$$

475 Because that $p, s = \Omega\left(\sqrt{\frac{\log(n/\delta)}{n}}\right)$, we further have,

$$n(p+s)/4 \leq D_i \leq 3n(p+s)/4.$$

476 Applying a union bound over $i \in [n]$ conclude the proof. \square

477 **Lemma G.4.** Suppose that $\delta > 0$ and $n \geq 8 \frac{p+s}{(p-s)^2} \log(4/\delta)$. Then with probability at least $1 - \delta$,

$$\frac{1}{2} \frac{p-s}{p+s} |y_i| \leq |\tilde{y}_i| \leq \frac{3}{2} \frac{p-s}{p+s} |y_i|.$$

478 *Proof of Lemma G.4.* By Hoeffding's inequality, with probability at least $1 - \delta/2$, we have

$$\left| \frac{1}{D_i} \sum_{k \in \mathcal{N}(i)} y_k - \frac{p-s}{p+s} y_i \right| \leq \sqrt{\frac{\log(4/\delta)}{2n(p+s)}}.$$

479 Therefore, as long as $n \geq 8 \frac{p+s}{(p-s)^2} \log(4/\delta)$, we have:

$$\frac{1}{2} \frac{p-s}{p+s} |y_i| \leq |\tilde{y}_i| \leq \frac{3}{2} \frac{p-s}{p+s} |y_i|.$$

480 This proves the result for the stability of sign of graph convoluted label. \square

481 **Lemma G.5.** Suppose that $\delta > 0$ and $d = \Omega(n^2(p+s)^2 \log(4n^2/\delta))$. Then with probability at least
482 $1 - \delta$,

$$\sigma_p^2 d / (4n(p+s)) \leq \|\tilde{\xi}_i\|_2^2 \leq 3\sigma_p^2 d / (4n(p+s)),$$

483 for all $i \in [n]$.

484 *Proof of Lemma G.5.* It is known that:

$$\|\tilde{\xi}_i\|_2^2 = \frac{1}{D_i^2} \sum_{j=1}^d \left(\sum_{k=1}^{D_i} \xi_{jk} \right)^2 = \frac{1}{D_i^2} \sum_{j=1}^d \sum_{k=1}^{D_i} \xi_{jk}^2 + \frac{1}{D_i^2} \sum_{j=1}^d \sum_{k \neq k'}^{D_i} \xi_{jk'} \xi_{jk}.$$

485 By Bernstein's inequality, with probability at least $1 - \delta/(2n)$ we have

$$\left| \sum_{j=1}^d \sum_{k=1}^{D_i} \xi_{jk}^2 - \sigma_p^2 d D_i \right| = O(\sigma_p^2 \cdot \sqrt{d D_i \log(4n/\delta)}).$$

486 Therefore, as long as $d = \Omega(\log(4n/\delta)/(n(p+s)))$, we have

$$3\sigma_p^2 d D_i / 4 \leq \sum_{j=1}^d \sum_{k=1}^{D_i} \xi_{jk}^2 \leq 5\sigma_p^2 d D_i / 4.$$

487 By Lemma G.3, we have,

$$2\sigma_p^2 d / (4n(p+s)) \leq \frac{1}{D_i^2} \sum_{j=1}^d \sum_{k=1}^{D_i} \xi_{jk}^2 \leq 6\sigma_p^2 d / (4n(p+s)).$$

488 Moreover, clearly $\langle \xi_k, \xi_{k'} \rangle$ has mean zero. For any k, k' with $k \neq k'$, by Bernstein's inequality, with
489 probability at least $1 - \delta / (2n^2)$ we have

$$|\langle \xi_k, \xi_{k'} \rangle| \leq 2\sigma_p^2 \cdot \sqrt{d \log(4n^2/\delta)}.$$

490 Applying a union bound we have that with probability at least $1 - \delta$,

$$|\langle \xi_k, \xi_{k'} \rangle| \leq 2\sigma_p^2 \cdot \sqrt{d \log(4n^2/\delta)}.$$

491 Therefore, as long as $d = \Omega(n^2(p+s)^2 \log(4n^2/\delta))$, we have

$$\sigma_p^2 d / (4n(p+s)) \leq \|\tilde{\xi}_i\|_2^2 \leq 3\sigma_p^2 d / (4n(p+s)).$$

492 **Remark G.6.** We compare the noise vector both before and after applying graph convolution. By
493 examining Lemma G.1 and Lemma G.5, we discover that the expectation of the ℓ_2 norm of noise
494 vector is reduced by a factor of $\sqrt{n(p+s)}/2$. This factor represents the square root of the expected
495 degree of the graph, indicating a significant change in the noise characteristics as a result of the
496 graph convolution process. We provide a demonstrative visualization in Figure 6.

497

□

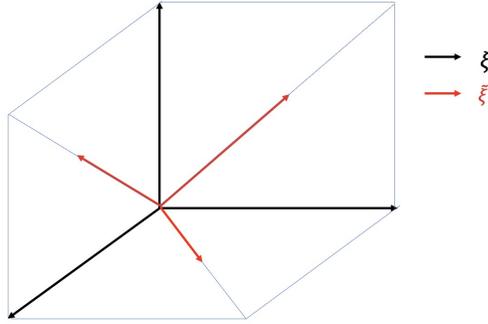


Figure 6: An illustrative example of noise vector before and after graph aggregation. In this example, we consider $d = 3$ and all degree are 1. The black vectors stand for noise vectors ξ before graph convolution. Each of them are orthogonal to each other. The red vectors represent noise vectors after graph convolution $\tilde{\xi}$. They are graph convoluted noise vectors of two original noise vectors. Note that the ℓ_2 norm between two kinds of vector follows $\|\tilde{\xi}\|_2 = \frac{\sqrt{2}}{2} \|\xi\|_2$. This plot demonstrates how graph convolution shrinks the ℓ_2 norm of noise vectors.

498 **Lemma G.7.** Suppose that $d = \Omega(n(p+s) \log(nm/\delta))$, $m = \Omega(\log(1/\delta))$. Then with probability
499 at least $1 - \delta$,

$$|\langle \mathbf{w}_{j,r}^{(0)}, \tilde{\xi}_i \rangle| \leq 4\sqrt{\log(8mn/\delta)} \cdot \sigma_0 \sigma_p \sqrt{d/(n(p+s))},$$

$$\sigma_0 \sigma_p \sqrt{d/(n(p+s))} / 4 \leq \max_{r \in [m]} j \cdot \langle \mathbf{w}_{j,r}^{(0)}, \tilde{\xi}_i \rangle \leq 2\sqrt{\log(8mn/\delta)} \cdot \sigma_0 \sigma_p \sqrt{d/(n(p+s))},$$

500 for all $j \in \{\pm 1\}$ and $i \in [n]$.

501 *Proof of Lemma G.7.* According to the fact that the weight $\mathbf{w}_{j,r}(0)$ and noise vector ξ are sampled
502 from Gaussian distribution, we know that $\langle \mathbf{w}_{j,r}^{(0)}, \tilde{\xi}_i \rangle$ is also Gaussian. By Lemma G.5, with probability
503 at least $1 - \delta/4$, we have that

$$\sigma_p \sqrt{d/(n(p+s))} / \sqrt{2} \leq \|\tilde{\xi}_i\|_2 \leq \sqrt{3/2} \cdot \sigma_p \sqrt{d/(n(p+s))}$$

504 holds for all $i \in [n]$. Therefore, applying the concentration bound for Gaussian variable, we obtain
 505 that

$$|\langle \mathbf{w}_{j,r}^{(0)}, \tilde{\boldsymbol{\xi}}_i \rangle| \leq 4\sqrt{\log(8mn/\delta)} \cdot \sigma_0 \sigma_p \sqrt{d/(n(p+s))}.$$

506 Next we finish the argument for the lower bound of maximum through the follow expression:

$$\begin{aligned} P(\max \langle \mathbf{w}_{j,r}^{(0)}, \tilde{\boldsymbol{\xi}}_i \rangle \geq \sigma_0 \sigma_p \sqrt{d/(n(p+s))}/4) &= 1 - P(\max \langle \mathbf{w}_{j,r}^{(0)}, \tilde{\boldsymbol{\xi}}_i \rangle < \sigma_0 \sigma_p \sqrt{d/(n(p+s))}/4) \\ &= 1 - P(\max \langle \mathbf{w}_{j,r}^{(0)}, \tilde{\boldsymbol{\xi}}_i \rangle < \sigma_0 \sigma_p \sqrt{d/(n(p+s))}/4)^{2m} \\ &\geq 1 - \delta/4. \end{aligned}$$

507 Together with Lemma G.5, we finally obtain that

$$\sigma_0 \sigma_p \sqrt{d/(n(p+s))}/4 \leq \max_{r \in [m]} j \cdot \langle \mathbf{w}_{j,r}^{(0)}, \tilde{\boldsymbol{\xi}}_i \rangle \leq 2\sqrt{\log(8mn/\delta)} \cdot \sigma_0 \sigma_p \sqrt{d/(n(p+s))}.$$

508

□

509 H General Lemmas for Iterative Coefficient Analysis

510 In this section, we deliver lemmas that delineate the iterative behavior of coefficients under gradient
 511 descent. We commence with proving the coefficient update rules as stated in Lemma E.1 in Section
 512 H.1. Subsequently, we establish the scale of training dynamics in Section H.2.

513 H.1 Coefficient update rule

514 **Lemma H.1** (Restatement of Lemma E.1). *The coefficients $\gamma_{j,r}^{(t)}$, $\bar{\rho}_{j,r,i}^{(t)}$, $\rho_{j,r,i}^{(t)}$ defined in Eq. (5) satisfy*
 515 *the following iterative equations:*

$$\begin{aligned} \gamma_{j,r}^{(0)}, \bar{\rho}_{j,r,i}^{(0)}, \rho_{j,r,i}^{(0)} &= 0, \\ \gamma_{j,r}^{(t+1)} &= \gamma_{j,r}^{(t)} - \frac{\eta}{nm} \cdot \sum_{i=1}^n \ell'_i{}^{(t)} \sigma'(\langle \mathbf{w}_{j,r}^{(t)}, \tilde{y}_i \boldsymbol{\mu} \rangle) y_i \tilde{y}_i \|\boldsymbol{\mu}\|_2^2, \\ \bar{\rho}_{j,r,i}^{(t+1)} &= \bar{\rho}_{j,r,i}^{(t)} - \frac{\eta}{nm} \cdot \sum_{k \in \mathcal{N}(i)} D_k^{-1} \cdot \ell'_k{}^{(t)} \cdot \sigma'(\langle \mathbf{w}_{j,r}^{(t)}, \tilde{\boldsymbol{\xi}}_k \rangle) \cdot \|\boldsymbol{\xi}_i\|_2^2 \cdot \mathbb{1}(y_k = j), \\ \rho_{j,r,i}^{(t+1)} &= \rho_{j,r,i}^{(t)} - \frac{\eta}{nm} \cdot \sum_{k \in \mathcal{N}(i)} D_k^{-1} \cdot \ell'_k{}^{(t)} \cdot \sigma'(\langle \mathbf{w}_{j,r}^{(t)}, \tilde{\boldsymbol{\xi}}_k \rangle) \cdot \|\boldsymbol{\xi}_i\|_2^2 \cdot \mathbb{1}(y_k = -j), \end{aligned}$$

516 for all $r \in [m]$, $j \in \{\pm 1\}$ and $i \in [n]$.

517 **Remark H.2.** *This lemma serves as a foundational element in our analysis of dynamics. Initially, the*
 518 *study of neural network dynamics under gradient descent required us to monitor the fluctuations in*
 519 *weights. However, this Lemma enables us to observe these dynamics through a new lens, focusing on*
 520 *two distinct aspects: signal learning and noise memorization. These are represented by the variables*
 521 *$\gamma_{j,r}^{(t)}$ and $\rho_{j,r,i}^{(t)}$, respectively. Furthermore, the selection of our data model was a conscious decision,*
 522 *designed to clearly separate the signal learning from the noise memorization aspects of learning.*
 523 *By maintaining a clear distinction between signal and noise, we can conduct a precise analysis of*
 524 *how each model learns the signal and memorizes the noise. This approach not only simplifies our*
 525 *understanding but also enhances our ability to dissect the underlying mechanisms of learning.*

526 *Proof of Lemma H.1.* Basically, the iteration of coefficients is derived based on gradient descent
 527 rule (4) and weight decomposition (5). We first consider $\hat{\gamma}_{j,r}^{(0)}, \hat{\rho}_{j,r,i}^{(0)} = 0$ and

$$\begin{aligned} \hat{\gamma}_{j,r}^{(t+1)} &= \hat{\gamma}_{j,r}^{(t)} - \frac{\eta}{nm} \cdot \sum_{i=1}^n \ell'_i{}^{(t)} \sigma'(\langle \mathbf{w}_{j,r}^{(t)}, \tilde{y}_i \boldsymbol{\mu} \rangle) y_i \tilde{y}_i \|\boldsymbol{\mu}\|_2^2, \\ \hat{\rho}_{j,r,i}^{(t+1)} &= \hat{\rho}_{j,r,i}^{(t)} - \frac{\eta}{nm} \cdot \sum_{k \in \mathcal{N}(i)} D_k^{-1} \cdot \ell'_k{}^{(t)} \cdot \sigma'(\langle \mathbf{w}_{j,r}^{(t)}, \tilde{\boldsymbol{\xi}}_k \rangle) \cdot \|\boldsymbol{\xi}_i\|_2^2 \cdot y_k, \end{aligned}$$

528 Taking above equations into Equation (4), we can obtain that

$$\mathbf{w}_{j,r}^{(t)} = \mathbf{w}_{j,r}^{(0)} + j \cdot \hat{\gamma}_{j,r}^{(t)} \cdot \|\boldsymbol{\mu}\|_2^{-2} \cdot \boldsymbol{\mu} + \sum_{i=1}^n \hat{\rho}_{j,r,i}^{(t)} \|\boldsymbol{\xi}_i\|_2^{-2} \cdot \boldsymbol{\xi}_i.$$

529 This result verifies that the iterative update of the coefficients is directly driven by the gradient
 530 descent update process. Furthermore, the uniqueness of the decomposition leads us to the precise
 531 relationships $\gamma_{j,r}^{(t)} = \hat{\gamma}_{j,r}^{(t)}$ and $\rho_{j,r,i}^{(t)} = \hat{\rho}_{j,r,i}^{(t)}$. Next, we examine the stability of the sign associated
 532 with noise memorization by employing the following telescopic analysis. This method allows us to
 533 investigate the continuity and consistency of the noise memorization process, providing insights into
 534 how the system behaves over successive iterations.

$$\rho_{j,r,i}^{(t)} = - \sum_{s=0}^{t-1} \sum_{k \in \mathcal{N}(i)} D_k^{-1} \frac{\eta}{nm} \cdot \ell_k^{(s)} \cdot \sigma'(\langle \mathbf{w}_{j,r}^{(s)}, \tilde{\boldsymbol{\xi}}_k \rangle) \cdot \|\boldsymbol{\xi}_i\|_2^2 \cdot j y_k.$$

535 Recall the sign of loss derivative is given by the definition of the cross-entropy loss, namely, $\ell_i^{(t)} < 0$.
 536 Therefore,

$$\bar{\rho}_{j,r,i}^{(t)} = - \sum_{s=0}^{t-1} \frac{\eta}{nm} \cdot \sum_{k \in \mathcal{N}(i)} D_k^{-1} \cdot \ell_k^{(s)} \cdot \sigma'(\langle \mathbf{w}_{j,r}^{(s)}, \tilde{\boldsymbol{\xi}}_k \rangle) \cdot \|\boldsymbol{\xi}_i\|_2^2 \cdot \mathbb{1}(y_k = j), \quad (13)$$

$$\rho_{j,r,i}^{(t)} = - \sum_{s=0}^{t-1} \frac{\eta}{nm} \cdot \sum_{k \in \mathcal{N}(i)} D_k^{-1} \cdot \ell_k^{(s)} \cdot \sigma'(\langle \mathbf{w}_{j,r}^{(s)}, \tilde{\boldsymbol{\xi}}_k \rangle) \cdot \|\boldsymbol{\xi}_i\|_2^2 \cdot \mathbb{1}(y_k = -j). \quad (14)$$

537 Writing out the iterative versions of (13) and (14) completes the proof. \square

538 **Remark H.3.** *The proof strategy follows the study of feature learning in CNN as described in [15].*
 539 *However, compared to CNNs, the decomposition of weights in GNN is notably more intricate. This*
 540 *complexity is particularly evident in the dynamics of noise memorization, as represented by Equations*
 541 *(13) and (14). The reason for this increased complexity lies in the additional graph convolution*
 542 *operations within GNNs. These operations introduce new interaction and dependencies, making the*
 543 *analysis of weight dynamics more challenging and nuanced.*

544 H.2 Scale of training dynamics

545 Our proof hinges on a meticulous evaluation of the coefficient values $\gamma_{j,r}^{(t)}$, $\bar{\rho}_{j,r,i}^{(t)}$, $\rho_{j,r,i}^{(t)}$ throughout the
 546 entire training process. In order to facilitate a more thorough analysis, we first establish the following
 547 bounds for these coefficients, which are maintained consistently throughout the training period.

548 Consider training the Graph Neural Network (GNN) for an extended period up to T^* . We aim to
 549 investigate the scale of noise memorization in relation to signal learning.

550 Let $T^* = \eta^{-1} \text{poly}(\epsilon^{-1}, \|\boldsymbol{\mu}\|_2^{-1}, d^{-1} \sigma_p^{-2}, \sigma_0^{-1}, n, m, d)$ be the maximum admissible iterations. De-
 551 note $\alpha = 4 \log(T^*)$. In preparation for an in-depth analysis, we enumerate the necessary conditions
 552 that must be satisfied. These conditions, which are essential for the subsequent examination, are also
 553 detailed in Condition 3.1:

$$\eta = O\left(\min\{nm/(q\sigma_p^2 d), nm/(q2^{q+2}\alpha^{q-2}\sigma_p^2 d), nm/(q2^{q+2}\alpha^{q-2}\|\boldsymbol{\mu}\|_2^2)\}\right), \quad (15)$$

$$\sigma_0 \leq [16\sqrt{\log(8mn/\delta)}]^{-1} \min\left\{\Xi^{-1}\|\boldsymbol{\mu}\|_2^{-1}, (\sigma_p \sqrt{d/(n(p+s))})^{-1}\right\}, \quad (16)$$

$$d \geq 1024 \log(4n^2/\delta) \alpha^2 n^2. \quad (17)$$

554 Denote $\beta = 2 \max_{i,j,r} \{|\langle \mathbf{w}_{j,r}^{(0)}, \tilde{y}_i \cdot \boldsymbol{\mu} \rangle|, |\langle \mathbf{w}_{j,r}^{(0)}, \tilde{\boldsymbol{\xi}}_i \rangle|\}$, it is straightforward to show the following
 555 inequality:

$$4 \max\left\{\beta, 8n \sqrt{\frac{\log(4n^2/\delta)}{d}} \alpha\right\} \leq 1. \quad (18)$$

556 First, by Lemma G.4 with probability at least $1 - \delta$, we can upper bound β by $4\sqrt{\log(8mn/\delta)} \cdot \sigma_0 \cdot$
557 $\max\{\Xi\|\boldsymbol{\mu}\|_2, \sigma_p\sqrt{d/(n(p+s))}\}$. Combined with the condition (16), we can bound β by 1. Second,
558 it is easy to check that $8n\sqrt{\frac{\log(4n^2/\delta)}{d}}\alpha \leq 1$ by inequality (17).

559 Having established the values of α and β at hand, we are now in a position to assert that the following
560 proposition holds for the entire duration of the training process, specifically for $0 \leq t \leq T^*$.

561 **Proposition H.4.** *Under Condition 3.1, for $0 \leq t \leq T^*$, where $T^* =$
562 $\eta^{-1}\text{poly}(\epsilon^{-1}, \|\boldsymbol{\mu}\|_2^{-1}, d^{-1}\sigma_p^{-2}, \sigma_0^{-1}, n, m, d)$, we have that*

$$0 \leq \gamma_{j,r}^{(t)}, \bar{\rho}_{j,r,i}^{(t)} \leq \alpha, \quad (19)$$

$$0 \geq \underline{\rho}_{j,r,i}^{(t)} \geq -\alpha, \quad (20)$$

563 for all $r \in [m]$, $j \in \{\pm 1\}$ and $i \in [n]$, where $\alpha = 4\log(T^*)$.

564 To establish Proposition H.4, we will employ an inductive approach. Before proceeding with the
565 proof, we need to introduce several technical lemmas that are fundamental to our argument.

566 We note that although the setting is slightly different from the case in [15]. With the same analysis,
567 we can obtain the following result.

568 **Lemma H.5** ([15]). *For any $t \geq 0$, it holds that $\langle \mathbf{w}_{j,r}^{(t)} - \mathbf{w}_{j,r}^{(0)}, \boldsymbol{\mu} \rangle = j \cdot \gamma_{j,r}^{(t)}$ for all $r \in [m]$, $j \in \{\pm 1\}$.*

569 In the subsequent three lemmas, our proof strategy is guided by the approach found in [15]. However,
570 we extend this methodology by providing a fine-grained analysis that takes into account the additional
571 complexity introduced by the graph convolution operation.

572 **Lemma H.6.** *Under Condition 3.1, suppose (19) and (20) hold at iteration t . Then*

$$\hat{\rho}_{j,r,i}^{(t)} - 8n\sqrt{\frac{\log(4n^2/\delta)}{d}}\alpha \leq \langle \mathbf{w}_{j,r}^{(t)} - \mathbf{w}_{j,r}^{(0)}, \tilde{\boldsymbol{\xi}}_i \rangle \leq \hat{\rho}_{j,r,i}^{(t)} + 8n\sqrt{\frac{\log(4n^2/\delta)}{d}}\alpha,$$

573 where $\hat{\rho}_{j,r,i} \triangleq \sum_{k \in \mathcal{N}(i)} D_i^{-1} \sum_{i' \neq k} \rho_{j,r,i'}^{(t)}$, for all $r \in [m]$, $j \in \{\pm 1\}$ and $i \in [n]$.

574 **Remark H.7.** *Lemma H.6 asserts that the inner product between the updated weight and the graph
575 convolution operation closely approximates the graph-convoluted noise memorization.*

576 *Proof of Lemma H.6.* It is known that,

$$\begin{aligned} \langle \mathbf{w}_{j,r}^{(t)} - \mathbf{w}_{j,r}^{(0)}, \tilde{\boldsymbol{\xi}}_i \rangle &= \sum_{i'=1}^n \bar{\rho}_{j,r,i'}^{(t)} \|\boldsymbol{\xi}_{i'}\|_2^{-2} \cdot \langle \boldsymbol{\xi}_{i'}, \tilde{\boldsymbol{\xi}}_i \rangle + \sum_{i'=1}^n \underline{\rho}_{j,r,i'}^{(t)} \|\boldsymbol{\xi}_{i'}\|_2^{-2} \cdot \langle \boldsymbol{\xi}_{i'}, \tilde{\boldsymbol{\xi}}_i \rangle \\ &= \sum_{i'=1}^n \sum_{k \in \mathcal{N}(i)} D_i^{-1} \bar{\rho}_{j,r,i'}^{(t)} \|\boldsymbol{\xi}_{i'}\|_2^{-2} \cdot \langle \boldsymbol{\xi}_{i'}, \boldsymbol{\xi}_k \rangle + \sum_{i'=1}^n \sum_{k \in \mathcal{N}(i)} D_i^{-1} \underline{\rho}_{j,r,i'}^{(t)} \|\boldsymbol{\xi}_{i'}\|_2^{-2} \cdot \langle \boldsymbol{\xi}_{i'}, \boldsymbol{\xi}_k \rangle \\ &\leq 4\sqrt{\frac{\log(4n^2/\delta)}{d}} \sum_{k \in \mathcal{N}(i)} D_i^{-1} \sum_{i' \neq k} |\bar{\rho}_{j,r,i'}^{(t)}| + 4\sqrt{\frac{\log(4n^2/\delta)}{d}} \sum_{k \in \mathcal{N}(i)} D_i^{-1} \sum_{i' \neq k} |\underline{\rho}_{j,r,i'}^{(t)}| \\ &\quad + \sum_{k \in \mathcal{N}(i)} D_i^{-1} \sum_{i' \neq k} \bar{\rho}_{j,r,i'}^{(t)} + \sum_{k \in \mathcal{N}(i)} D_i^{-1} \sum_{i' \neq k} \underline{\rho}_{j,r,i'}^{(t)} \\ &\leq \hat{\rho}_{j,r,i}^{(t)} + 8n\sqrt{\frac{\log(4n^2/\delta)}{d}}\alpha, \end{aligned}$$

577 where we define $\hat{\rho}_{j,r,i} \triangleq \sum_{k \in \mathcal{N}(i)} D_i^{-1} \sum_{i' \neq k} \rho_{j,r,i'}^{(t)}$ the second inequality is by Lemma G.1 and the
578 last inequality is by $|\bar{\rho}_{j,r,i'}^{(t)}|, |\underline{\rho}_{j,r,i'}^{(t)}| \leq \alpha$ in (19).

579 Similarly, we can show that:

$$\begin{aligned}
\langle \mathbf{w}_{j,r}^{(t)} - \mathbf{w}_{j,r}^{(0)}, \tilde{\boldsymbol{\xi}}_i \rangle &= \sum_{i'=1}^n \bar{\rho}_{j,r,i'}^{(t)} \|\boldsymbol{\xi}_{i'}\|_2^{-2} \cdot \langle \boldsymbol{\xi}_{i'}, \tilde{\boldsymbol{\xi}}_i \rangle + \sum_{i'=1}^n \underline{\rho}_{j,r,i'}^{(t)} \|\boldsymbol{\xi}_{i'}\|_2^{-2} \cdot \langle \boldsymbol{\xi}_{i'}, \tilde{\boldsymbol{\xi}}_i \rangle \\
&= \sum_{i'=1}^n \sum_{k \in \mathcal{N}(i)} D_i^{-1} \bar{\rho}_{j,r,i'}^{(t)} \|\boldsymbol{\xi}_{i'}\|_2^{-2} \cdot \langle \boldsymbol{\xi}_{i'}, \boldsymbol{\xi}_k \rangle + \sum_{i'=1}^n \sum_{k \in \mathcal{N}(i)} D_i^{-1} \underline{\rho}_{j,r,i'}^{(t)} \|\boldsymbol{\xi}_{i'}\|_2^{-2} \cdot \langle \boldsymbol{\xi}_{i'}, \boldsymbol{\xi}_k \rangle \\
&\geq -4\sqrt{\frac{\log(4n^2/\delta)}{d}} \sum_{k \in \mathcal{N}(i)} D_i^{-1} \sum_{i' \neq k} |\bar{\rho}_{j,r,i'}^{(t)}| - 4\sqrt{\frac{\log(4n^2/\delta)}{d}} \sum_{k \in \mathcal{N}(i)} D_i^{-1} \sum_{i' \neq k} |\underline{\rho}_{j,r,i'}^{(t)}| \\
&\quad + \sum_{k \in \mathcal{N}(i)} D_i^{-1} \sum_{i' \neq k} \bar{\rho}_{j,r,i'}^{(t)} + \sum_{k \in \mathcal{N}(i)} D_i^{-1} \sum_{i' \neq k} \underline{\rho}_{j,r,i'}^{(t)} \\
&\geq \hat{\rho}_{j,r,i}^{(t)} - 8n\sqrt{\frac{\log(4n^2/\delta)}{d}} \alpha,
\end{aligned}$$

580 where the first inequality is by Lemma G.1 and the second inequality is by $|\bar{\rho}_{j,r,i'}^{(t)}|, |\underline{\rho}_{j,r,i'}^{(t)}| \leq \alpha$ in
581 (19), which completes the proof. \square

582 **Lemma H.8.** *Under Condition 3.1, suppose (19) and (20) hold at iteration t . Then*

$$\begin{aligned}
\langle \mathbf{w}_{j,r}^{(t)}, \tilde{y}_i \boldsymbol{\mu} \rangle &\leq \langle \mathbf{w}_{j,r}^{(0)}, \tilde{y}_i \boldsymbol{\mu} \rangle, \\
\langle \mathbf{w}_{j,r}^{(t)}, \tilde{\boldsymbol{\xi}}_i \rangle &\leq \langle \mathbf{w}_{j,r}^{(0)}, \tilde{\boldsymbol{\xi}}_i \rangle + 8n\sqrt{\frac{\log(4n^2/\delta)}{d}} \alpha,
\end{aligned}$$

583 for all $r \in [m]$ and $j \neq y_i$. If $\max\{\gamma_{j,r}^{(t)}, \rho_{j,r,i}^{(t)}\} = O(1)$, we further have that $F_j(\mathbf{W}_j^{(t)}, \tilde{\mathbf{x}}_i) = O(1)$.

584 **Remark H.9.** *Lemma H.8 further establishes that the update in the direction of $\tilde{\boldsymbol{\xi}}$ can be constrained*
585 *within specific bounds when $j \neq y_i$. As a result, the output function remains controlled and does not*
586 *exceed a constant order.*

587 *Proof of Lemma H.8.* For $j \neq y_i$, we have that

$$\langle \mathbf{w}_{j,r}^{(t)}, \tilde{y}_i \boldsymbol{\mu} \rangle = \langle \mathbf{w}_{j,r}^{(0)}, \tilde{y}_i \boldsymbol{\mu} \rangle + \tilde{y}_i \cdot j \cdot \gamma_{j,r}^{(t)} \leq \langle \mathbf{w}_{j,r}^{(0)}, \tilde{y}_i \boldsymbol{\mu} \rangle, \quad (21)$$

588 where the inequality is by $\gamma_{j,r}^{(t)} \geq 0$ and Lemma G.4 stating that $\text{sign}(y_i) = \text{sign}(\tilde{y}_i)$ with a high
589 probability. In addition, we have

$$\begin{aligned}
\langle \mathbf{w}_{j,r}^{(t)}, \tilde{\boldsymbol{\xi}}_i \rangle &= \langle \mathbf{w}_{j,r}^{(0)}, \tilde{\boldsymbol{\xi}}_i \rangle + \sum_{k \in \mathcal{N}(i)} D_i^{-1} \sum_{i'=1}^n \rho_{j,r,i'} \langle \boldsymbol{\xi}_k, \boldsymbol{\xi}_{i'} \rangle \|\boldsymbol{\xi}_{i'}\|_2^{-2} \\
&\leq \langle \mathbf{w}_{j,r}^{(0)}, \tilde{\boldsymbol{\xi}}_i \rangle + D_i^{-1} \left(\sum_{y_k \neq j} \rho_{j,r,i}^{(t)} + \sum_{y_k = j} \bar{\rho}_{j,r,i}^{(t)} \right) + 8n\sqrt{\frac{\log(4n^2/\delta)}{d}} \alpha \\
&\leq \langle \mathbf{w}_{j,r}^{(0)}, \tilde{\boldsymbol{\xi}}_i \rangle + 8n\sqrt{\frac{\log(4n^2/\delta)}{d}} \alpha,
\end{aligned} \quad (22)$$

590 where the first inequality is by Lemma H.6 and the second inequality is due to $\hat{\rho}_{j,r,i}^{(t)} \leq 0$ based on
 591 Lemma G.4. Then we can get that

$$\begin{aligned}
 F_j(\mathbf{W}_j^{(t)}, \tilde{\mathbf{x}}_i) &= \frac{1}{m} \sum_{r=1}^m [\sigma(\langle \mathbf{w}_{j,r}^{(t)}, \tilde{\mathbf{y}}_i \cdot \boldsymbol{\mu} \rangle) + \sigma(\langle \mathbf{w}_{j,r}^{(t)}, \tilde{\boldsymbol{\xi}}_i \rangle)] \\
 &= \frac{1}{m} \sum_{r=1}^m [\sigma(\langle \mathbf{w}_{j,r}^{(t)}, \tilde{\mathbf{y}}_i \cdot \boldsymbol{\mu} \rangle) + \sigma(\langle \mathbf{w}_{j,r}^{(t)}, D_i^{-1} \sum_{k \in \mathcal{N}(i)} \boldsymbol{\xi}_k \rangle)] \\
 &= \frac{1}{m} \sum_{r=1}^m [\sigma(\langle \mathbf{w}_{j,r}^{(0)}, \tilde{\mathbf{y}}_i \cdot \boldsymbol{\mu} \rangle) + \sigma(\langle \mathbf{w}_{j,r}^{(0)}, \tilde{\boldsymbol{\xi}}_i \rangle + \langle \mathbf{w}_{j,r}^{(t)} - \mathbf{w}_{j,r}^{(0)}, D_i^{-1} \sum_{k \in \mathcal{N}(i)} \boldsymbol{\xi}_k \rangle)] \\
 &\leq \frac{1}{m} \sum_{r=1}^m [\sigma(\langle \mathbf{w}_{j,r}^{(0)}, \tilde{\mathbf{y}}_i \cdot \boldsymbol{\mu} \rangle) + \sigma(\langle \mathbf{w}_{j,r}^{(0)}, \tilde{\boldsymbol{\xi}}_i \rangle) + 8n \sqrt{\frac{\log(4n^2/\delta)}{d}} \alpha + \hat{\rho}_{j,r,i}^{(t)}] \\
 &\leq 2^{q+1} \max_{j,r,i} \left\{ |\langle \mathbf{w}_{j,r}^{(0)}, \tilde{\mathbf{y}}_i \cdot \boldsymbol{\mu} \rangle|, |\langle \mathbf{w}_{j,r}^{(0)}, \tilde{\boldsymbol{\xi}}_i \rangle|, 8n \sqrt{\frac{\log(4n^2/\delta)}{d}} \alpha \right\}^q \\
 &\leq 1,
 \end{aligned}$$

592 where the first inequality is by (21), (22) and the second inequality is by (18) and $\max\{\gamma_{j,r}^{(t)}, \rho_{j,r,i}^{(t)}\} =$
 593 $O(1)$. \square

594 **Lemma H.10.** *Under Condition 3.1, suppose (19) and (20) hold at iteration t . Then*

$$\begin{aligned}
 \langle \mathbf{w}_{j,r}^{(t)}, \tilde{\mathbf{y}}_i \boldsymbol{\mu} \rangle &= \langle \mathbf{w}_{j,r}^{(0)}, \tilde{\mathbf{y}}_i \boldsymbol{\mu} \rangle + \gamma_{j,r}^{(t)}, \\
 \langle \mathbf{w}_{j,r}^{(t)}, \tilde{\boldsymbol{\xi}}_i \rangle &\leq \langle \mathbf{w}_{j,r}^{(0)}, \tilde{\boldsymbol{\xi}}_i \rangle + \hat{\rho}_{j,r,i}^{(t)} + 8n \sqrt{\frac{\log(4n^2/\delta)}{d}} \alpha
 \end{aligned}$$

595 for all $r \in [m]$, $j = y_i$ and $i \in [n]$. If $\max\{\gamma_{j,r}^{(t)}, \rho_{j,r,i}^{(t)}\} = O(1)$, we further have that
 596 $F_j(\mathbf{W}_j^{(t)}, \tilde{\mathbf{x}}_i) = O(1)$.

597 **Remark H.11.** *Lemma H.10 further establishes that the update in the direction of $\boldsymbol{\mu}$ and $\tilde{\boldsymbol{\xi}}$ can be*
 598 *constrained within specific bounds when $j = y_i$. As a result, the output function remains controlled*
 599 *and does not exceed a constant order with an additional condition.*

600 *Proof of Lemma H.10.* For $j = y_i$, we have that

$$\langle \mathbf{w}_{j,r}^{(t)}, \tilde{\mathbf{y}}_i \boldsymbol{\mu} \rangle = \langle \mathbf{w}_{j,r}^{(0)}, \tilde{\mathbf{y}}_i \boldsymbol{\mu} \rangle + \gamma_{j,r}^{(t)}, \tag{23}$$

601 where the equation is by Lemma H.5. We also have that

$$\langle \mathbf{w}_{j,r}^{(t)}, \tilde{\boldsymbol{\xi}}_i \rangle \leq \langle \mathbf{w}_{j,r}^{(0)}, \tilde{\boldsymbol{\xi}}_i \rangle + \hat{\rho}_{j,r,i}^{(t)} + 8n \sqrt{\frac{\log(4n^2/\delta)}{d}} \alpha, \tag{24}$$

602 where the inequality is by Lemma H.6. If $\max\{\gamma_{j,r}^{(t)}, \rho_{j,r,i}^{(t)}\} = O(1)$, we have following bound

$$\begin{aligned}
 F_j(\mathbf{W}_j^{(t)}, \tilde{\mathbf{x}}_i) &= \frac{1}{m} \sum_{r=1}^m [\sigma(\langle \mathbf{w}_{j,r}^{(t)}, \tilde{\mathbf{y}}_i \cdot \boldsymbol{\mu} \rangle) + \sigma(\langle \mathbf{w}_{j,r}^{(t)}, \tilde{\boldsymbol{\xi}}_i \rangle)] \\
 &\leq 2 \cdot 3^q \max_{j,r,i} \left\{ \gamma_{j,r}^{(t)}, |\hat{\rho}_{j,r,i}^{(t)}|, |\langle \mathbf{w}_{j,r}^{(0)}, \tilde{\mathbf{y}}_i \cdot \boldsymbol{\mu} \rangle|, |\langle \mathbf{w}_{j,r}^{(0)}, \tilde{\boldsymbol{\xi}}_i \rangle|, 8n \sqrt{\frac{\log(4n^2/\delta)}{d}} \alpha \right\}^q \\
 &= O(1),
 \end{aligned}$$

603 where $\hat{\rho}_{j,r,i}^{(t)} = \frac{1}{D_i} \sum_{k \in \mathcal{N}(i)} \bar{\rho}_{j,r,k}^{(t)} \mathbb{1}(y_k = j) + \bar{\rho}_{j,r,k}^{(t)} \mathbb{1}(y_k \neq j)$, the first inequality is by (23), (24).
 604 Then the second inequality is by (18) where $\beta = 2 \max_{i,j,r} \{|\langle \mathbf{w}_{j,r}^{(0)}, \tilde{\mathbf{y}}_i \cdot \boldsymbol{\mu} \rangle|, |\langle \mathbf{w}_{j,r}^{(0)}, \tilde{\boldsymbol{\xi}}_i \rangle|\} \leq 1$ and
 605 condition that $\max\{\gamma_{j,r}^{(t)}, \rho_{j,r,i}^{(t)}\} = O(1)$. \square

606 Equipped with Lemmas H.5 - H.10, we are now prepared to prove Proposition H.4. These lemmas
607 provide the foundational building blocks and insights necessary for our proof, setting the stage for a
608 rigorous and comprehensive demonstration of the proposition

609 *Proof of Proposition H.4.* Following a similar approach to the proof found in [15], we employ an
610 induction method. This technique allows us to build our argument step by step, drawing on established
611 principles and extending them to our specific context, thereby providing a robust and systematic
612 demonstration.

613 At the initial time step $t = 0$, the outcome is clear since all coefficients are set to zero.

614 Next, we hypothesize that there exists a time \tilde{T} less than T^* during which Proposition H.4 holds true
615 for every moment within the range $0 \leq t \leq \tilde{T} - 1$. Our objective is to show that this proposition
616 remains valid at $t = \tilde{T}$.

617 We aim to validate that equation (20) is applicable at $t = \tilde{T}$, meaning that,

$$\rho_{j,r,i}^{(t)} \geq -\beta - 16n\sqrt{\frac{\log(4n^2/\delta)}{d}}\alpha,$$

618 for the given parameters. It's important to note that $\rho_{j,r,i}^{(t)} = 0$ when $j = y_i$. So we only need to
619 consider instances where $j \neq y_i$.

620 1) Under condition

$$\rho_{j,r,i}^{(\tilde{T}-1)} \leq -0.5\beta - 8n\sqrt{\frac{\log(4n^2/\delta)}{d}}\alpha,$$

621 Lemma H.6 leads us to the following relationships:

$$\langle \mathbf{w}_{j,r}^{(\tilde{T}-1)}, \tilde{y}_i \boldsymbol{\mu} \rangle \leq \rho_{j,r,i}^{(\tilde{T}-1)} + \langle \mathbf{w}_{j,r}^{(0)}, \tilde{y}_i \boldsymbol{\mu} \rangle + 8n\sqrt{\frac{\log(4n^2/\delta)}{d}}\alpha \leq 0,$$

622 and thus

$$\begin{aligned} \rho_{j,r,i}^{(\tilde{T})} &= \rho_{j,r,i}^{(\tilde{T}-1)} + \frac{\eta}{nm} \sum_k D_k^{-1} \cdot \ell_k'(\tilde{T}-1) \cdot \sigma'(\langle \mathbf{w}_{j,r}^{(\tilde{T}-1)}, \tilde{\boldsymbol{\xi}}_k \rangle) \cdot \mathbb{1}(y_k = -j) \|\boldsymbol{\xi}_i\|_2^2 \\ &= \rho_{j,r,i}^{(\tilde{T}-1)} \geq -\beta - 16n\sqrt{\frac{\log(4n^2/\delta)}{d}}\alpha, \end{aligned}$$

623 with the final inequality being supported by the induction hypothesis.

624 2) Given the condition $\rho_{j,r,i}^{(\tilde{T}-1)} \geq -0.5\beta - 8n\sqrt{\frac{\log(4n^2/\delta)}{d}}\alpha$, we can derive the following:

$$\begin{aligned} \rho_{j,r,i}^{(\tilde{T})} &= \rho_{j,r,i}^{(\tilde{T}-1)} + \frac{\eta}{nm} \cdot \sum_{k \in \mathcal{N}(i)} D_k^{-1} \ell_k'(\tilde{T}-1) \cdot \sigma'(\langle \mathbf{w}_{j,r}^{(\tilde{T}-1)}, \tilde{\boldsymbol{\xi}}_k \rangle) \cdot \mathbb{1}(y_k = -j) \|\boldsymbol{\xi}_i\|_2^2 \\ &\geq -0.5\beta - 8n\sqrt{\frac{\log(4n^2/\delta)}{d}}\alpha - O\left(\frac{\eta\sigma_p^2 d}{nm}\right) \sigma' \left(0.5\beta + 8n\sqrt{\frac{\log(4n^2/\delta)}{d}}\alpha\right) \\ &\geq -0.5\beta - 8n\sqrt{\frac{\log(4n^2/\delta)}{d}}\alpha - O\left(\frac{\eta q \sigma_p^2 d}{nm}\right) \left(0.5\beta + 8n\sqrt{\frac{\log(4n^2/\delta)}{d}}\alpha\right) \\ &\geq -\beta - 16n\sqrt{\frac{\log(4n^2/\delta)}{d}}\alpha, \end{aligned}$$

625 where we apply the inequalities $\ell_i'(\tilde{T}-1) \leq 1$ and $\|\boldsymbol{\xi}_i\|_2 = O(\sigma_p^2 d)$, and use the conditions $\eta =$

626 $O(nm/(q\sigma_p^2 d))$ and $0.5\beta + 8n\sqrt{\frac{\log(4n^2/\delta)}{d}}\alpha \leq 1$, as specified in (15).

627 Next, we aim to show that (19) is valid for $t = \tilde{T}$. We can express:

$$\begin{aligned} |\ell_i^{(t)}| &= \frac{1}{1 + \exp\{y_i \cdot [F_{+1}(\mathbf{W}_{+1}^{(t)}, \tilde{\mathbf{x}}_i) - F_{-1}(\mathbf{W}_{-1}^{(t)}, \tilde{\mathbf{x}}_i)]\}} \\ &\leq \exp\{-y_i \cdot [F_{+1}(\mathbf{W}_{+1}^{(t)}, \tilde{\mathbf{x}}_i) - F_{-1}(\mathbf{W}_{-1}^{(t)}, \tilde{\mathbf{x}}_i)]\} \\ &\leq \exp\{-F_{y_i}(\mathbf{W}_{y_i}^{(t)}, \tilde{\mathbf{x}}_i) + 1\}. \end{aligned} \tag{25}$$

628 with the last inequality being a result of Lemma H.8. Additionally, we recall the update rules for
 629 $\gamma_{j,r}^{(t+1)}$ and $\bar{\rho}_{j,r,i}^{(t+1)}$:

$$\begin{aligned}\gamma_{j,r}^{(t+1)} &= \gamma_{j,r}^{(t)} - \frac{\eta}{nm} \cdot \sum_{i=1}^n \ell_i^{(t)} \cdot \sigma'(\langle \mathbf{w}_{j,r}^{(t)}, \tilde{\mathbf{y}}_i \cdot \boldsymbol{\mu} \rangle) y_i \tilde{\mathbf{y}}_i \|\boldsymbol{\mu}\|_2^2, \\ \bar{\rho}_{j,r,i}^{(t+1)} &= \bar{\rho}_{j,r,i}^{(t)} - \frac{\eta}{nm} \cdot \sum_{k \in \mathcal{N}(i)} D_k^{-1} \ell_k^{(t)} \cdot \sigma'(\langle \mathbf{w}_{j,r}^{(t)}, \tilde{\boldsymbol{\xi}}_k \rangle) \cdot \mathbb{1}(y_k = j) \|\boldsymbol{\xi}_i\|_2^2.\end{aligned}$$

630 We define $t_{j,r,i}$ as the final moment $t < T^*$ when $\bar{\rho}_{j,r,i}^{(t)} \leq 0.5\alpha$.

631 We can express $\bar{\rho}_{j,r,i}^{(\tilde{T})}$ as follows:

$$\begin{aligned}\bar{\rho}_{j,r,i}^{(\tilde{T})} &= \bar{\rho}_{j,r,i}^{(t_{j,r,i})} - \underbrace{\frac{\eta}{nm} \cdot \sum_{k \in \mathcal{N}(i)} D_k^{-1} \cdot \ell_k^{(t_{j,r,i})} \cdot \sigma'(\langle \mathbf{w}_{j,r}^{(t_{j,r,i})}, \tilde{\boldsymbol{\xi}}_k \rangle) \cdot \mathbb{1}(y_k = j) \|\boldsymbol{\xi}_i\|_2^2}_{I_1} \\ &\quad - \underbrace{\sum_{t_{j,r,i} < t < T} \frac{\eta}{nm} \cdot \sum_{k \in \mathcal{N}(i)} D_k^{-1} \cdot \ell_k^{(t)} \cdot \sigma'(\langle \mathbf{w}_{j,r}^{(t)}, \tilde{\boldsymbol{\xi}}_k \rangle) \cdot \mathbb{1}(y_k = j) \|\boldsymbol{\xi}_i\|_2^2}_{I_2}.\end{aligned}\quad (26)$$

632 Next, we aim to establish an upper bound for I_1 :

$$\begin{aligned}|I_1| &\leq 2qn^{-1}m^{-1}\eta \left(\max_k \hat{\rho}_{j,r,k}^{(t_{j,r,i})} + 0.5\beta + 8n\sqrt{\frac{\log(4n^2/\delta)}{d}}\alpha \right)^{q-1} \sigma_p^2 d \\ &\leq q2^q n^{-1} m^{-1} \eta \alpha^{q-1} \sigma_p^2 d \leq 0.25\alpha,\end{aligned}$$

633 where we apply Lemmas H.6 and G.1 for the first inequality, utilize the conditions $\beta \leq$
 634 0.1α and $8n\sqrt{\frac{\log(4n^2/\delta)}{d}}\alpha \leq 0.1\alpha$ for the second inequality, and finally, the constraint $\eta \leq$
 635 $nm/(q2^{q+2}\alpha^{q-2}\sigma_p^2 d)$ for the last inequality.

636 Second, we bound I_2 . For $t_{j,r,i} < t < \tilde{T}$ and $y_k = j$, we can lower bound $\langle \mathbf{w}_{j,r}^{(t)}, \tilde{\boldsymbol{\xi}}_k \rangle$ as follows,

$$\begin{aligned}\langle \mathbf{w}_{j,r}^{(t)}, \tilde{\boldsymbol{\xi}}_k \rangle &\geq \langle \mathbf{w}_{j,r}^{(0)}, \tilde{\boldsymbol{\xi}}_k \rangle + \hat{\rho}_{j,r,k}^{(t)} - 8n\sqrt{\frac{\log(4n^2/\delta)}{d}}\alpha \\ &\geq -0.5\beta + \frac{1}{4}\frac{p-s}{p+s}\alpha - 8n\sqrt{\frac{\log(4n^2/\delta)}{d}}\alpha \\ &\geq 0.25\alpha,\end{aligned}$$

637 where the first inequality is by Lemma H.6, the second inequality is by $\hat{\rho}_{j,r,i}^{(t)} > \frac{1}{4}\frac{p-s}{p+s}\alpha$ and
 638 $\langle \mathbf{w}_{j,r}^{(0)}, \tilde{\boldsymbol{\xi}}_i \rangle \geq -0.5\beta$ due to the definition of $t_{j,r,i}$ and β , the last inequality is by $\beta \leq 0.1\alpha$ and
 639 $8n\sqrt{\frac{\log(4n^2/\delta)}{d}}\alpha \leq 0.1\alpha$. Similarly, for $t_{j,r,i} < t < \tilde{T}$ and $y_k = j$, we can also upper bound
 640 $\langle \mathbf{w}_{j,r}^{(t)}, \tilde{\boldsymbol{\xi}}_k \rangle$ as follows,

$$\begin{aligned}\langle \mathbf{w}_{j,r}^{(t)}, \tilde{\boldsymbol{\xi}}_k \rangle &\leq \langle \mathbf{w}_{j,r}^{(0)}, \tilde{\boldsymbol{\xi}}_k \rangle + \hat{\rho}_{j,r,k}^{(t)} + 8n\sqrt{\frac{\log(4n^2/\delta)}{d}}\alpha \\ &\leq 0.5\beta + \frac{3}{4}\frac{p-s}{p+s}\alpha + 8n\sqrt{\frac{\log(4n^2/\delta)}{d}}\alpha \\ &\leq 2\alpha,\end{aligned}$$

641 where the first inequality is by Lemma H.6, the second inequality is by induction hypothesis $\hat{\rho}_{j,r,i}^{(t)} \leq \alpha$,
 642 the last inequality is by $\beta \leq 0.1\alpha$ and $8n\sqrt{\frac{\log(4n^2/\delta)}{d}}\alpha \leq 0.1\alpha$.

643 Hence, we can derive the following expression for I_2 :

$$\begin{aligned}
|I_2| &\leq \sum_{t_{j,r,i} < t < \bar{T}} \frac{\eta}{nm} \cdot \sum_{k \in \mathcal{N}(i)} D_k^{-1} \exp(-\sigma(\langle \mathbf{w}_{j,r}^{(t)}, \tilde{\xi}_k \rangle) + 1) \cdot \sigma'(\langle \mathbf{w}_{j,r}^{(t)}, \tilde{\xi}_k \rangle) \cdot \mathbb{1}(y_k = j) \|\xi_i\|_2^2 \\
&\leq \frac{eq2^q \eta T^*}{n} \exp(-\alpha^q/4^q) \alpha^{q-1} \sigma_p^2 d \\
&\leq 0.25 T^* \exp(-\alpha^q/4^q) \alpha \\
&\leq 0.25 T^* \exp(-\log(T^*)^q) \alpha \\
&\leq 0.25 \alpha,
\end{aligned}$$

644 where we apply (25) for the first inequality, utilize Lemma G.1 for the second, employ the constraint
645 $\eta = O(nm/(q2^{q+2}\alpha^{q-2}\sigma_p^2d))$ in (15) for the third, and finally, the conditions $\alpha = 4 \log(T^*)$ and
646 $\log(T^*)^q \geq \log(T^*)$ for the subsequent inequalities. By incorporating the bounds of I_1 and I_2 into
647 (26), we conclude the proof for $\bar{\rho}$.

648 In a similar manner, we can establish that $\gamma_{j,r}^{(\bar{T})} \leq \alpha$ by using $\eta = O(nm/(q2^{q+2}\alpha^{q-2}\|\mu\|_2^2))$ in
649 (15). Thus, Proposition H.4 is valid for $t = \bar{T}$, completing the induction process. As a corollary to
650 Proposition H.4, we identify a crucial characteristic of the loss function during training within the
651 interval $0 \leq t \leq T^*$. This characteristic will play a vital role in the subsequent convergence analysis.

652 □

653 I Two Stage Dynamics Analysis

654 In this section, we employ a two-stage dynamics analysis to investigate the behavior of coefficient
655 iterations. During the first stage, the derivative of the loss function remains almost constant due to
656 the small weight initialization. In the second stage, the derivative of the loss function ceases to be
657 constant, necessitating an analysis that meticulously takes this into account.

658 I.1 First stage: feature learning versus noise memorization

659 **Lemma I.1** (Restatement of Lemma E.2). *Under the same conditions as Theorem 3.2, in particular*
660 *if we choose*

$$n \cdot \text{SNR}^q \cdot (n(p+s))^{q/2-1} \geq C \log(6/\sigma_0 \|\mu\|_2) 2^{2q+6} [4 \log(8mn/\delta)]^{(q-1)/2}, \quad (27)$$

661 where $C = O(1)$ is a positive constant, there exists time $T_1 = \frac{C \log(6/\sigma_0 \|\mu\|_2) 2^{q+1} m}{\eta \sigma_0^{q-2} \|\mu\|_2^q \Xi^q}$ such that

- 662 • $\max_r \gamma_{j,r}^{(T_1)} \geq 2$ for $j \in \{\pm 1\}$.
- 663 • $|\rho_{j,r,i}^{(t)}| \leq \sigma_0 \sigma_p \sqrt{d/(n(p+s))}/2$ for all $j \in \{\pm 1\}, r \in [m], i \in [n]$ and $0 \leq t \leq T_1$.

664 **Remark I.2.** *In this lemma, we establish that the rate of signal learning significantly outpaces that of*
665 *noise memorization within GNNs. After a specific number of iterations, the GNN is able to learn the*
666 *signal from the data at a constant or higher order, while only memorizing a smaller order of noise.*

667 *Proof of Lemma I.1.* Let us define

$$T_1^+ = \frac{nm\eta^{-1} \sigma_0^{2-q} \sigma_p^{-q} d^{-q/2} (n(p+s))^{(q-2)/2}}{2^{q+4} q [4 \log(8mn/\delta)]^{(q-2)/2}}. \quad (28)$$

668 We will begin by establishing the outcome related to noise memorization. Let $\Psi^{(t)}$ be the maximum
669 value over all j, r, i of $|\rho_{j,r,i}^{(t)}|$, that is, $\Psi^{(t)} = \max_{j,r,i} \{\bar{\rho}_{j,r,i}^{(t)}, -\underline{\rho}_{j,r,i}^{(t)}\}$. We will employ an inductive
670 argument to demonstrate that

$$\Psi^{(t)} \leq \sigma_0 \sigma_p \sqrt{d/(n(p+s))} \quad (29)$$

671 is valid for the entire range $0 \leq t \leq T_1^+$. By its very definition, it is evident that $\Psi^{(0)} = 0$. Assuming
672 that there exists a value $\tilde{T} \leq T_1^+$ for which equation (29) is satisfied for all $0 < t \leq \tilde{T} - 1$, we can
673 proceed as follows.

$$\begin{aligned}
\Psi^{(t+1)} &\leq \Psi^{(t)} + \frac{\eta}{nm} \sum_{k \in \mathcal{N}(i)} D_k^{-1} \cdot |\ell_k^{(t)}|. \\
&\sigma' \left(\langle \mathbf{w}_{j,r}^{(0)}, \tilde{\xi}_k \rangle + \sum_{i'=1}^n \Psi^{(t)} \cdot \frac{|\langle \xi_{i'}, \tilde{\xi}_k \rangle|}{\|\xi_{i'}\|_2^2} + \sum_{i'=1}^n \Psi^{(t)} \cdot \frac{|\langle \xi_{i'}, \tilde{\xi}_k \rangle|}{\|\xi_{i'}\|_2^2} \right) \cdot \|\xi_i\|_2^2 \\
&\leq \Psi^{(t)} + \frac{\eta}{nm} \cdot \sum_{k \in \mathcal{N}(i)} D_k^{-1} \sigma' \left(\langle \mathbf{w}_{j,r}^{(0)}, \tilde{\xi}_k \rangle + 2 \cdot \sum_{i'=1}^n \Psi^{(t)} \cdot \frac{|\langle \xi_{i'}, \tilde{\xi}_k \rangle|}{\|\xi_{i'}\|_2^2} \right) \cdot \|\xi_i\|_2^2 \\
&= \Psi^{(t)} + \frac{\eta}{nm} \cdot \sum_{k \in \mathcal{N}(i)} D_k^{-1} \cdot \\
&\sigma' \left(\langle \mathbf{w}_{j,r}^{(0)}, \tilde{\xi}_k \rangle + 2\Psi^{(t)} + 2 \cdot \sum_{i' \neq k'}^n \Psi^{(t)} \cdot D_k^{-1} \sum_{k' \in \mathcal{N}(k)} \frac{|\langle \xi_{i'}, \xi_{k'} \rangle|}{\|\xi_{i'}\|_2^2} \right) \cdot \|\xi_i\|_2^2 \\
&\leq \Psi^{(t)} + \frac{\eta q}{nm} \cdot \sum_{k \in \mathcal{N}(i)} D_k^{-1} \left[2 \cdot \sqrt{\log(8mn/\delta)} \cdot \sigma_0 \sigma_p \sqrt{d/(n(p+s))} \right. \\
&\quad \left. + \left(2 + \frac{4n\sigma_p^2 \cdot \sqrt{d \log(4n^2/\delta)}}{\sigma_p^2 d} \right) \cdot \Psi^{(t)} \right]^{q-1} \cdot 2\sigma_p^2 d \\
&\leq \Psi^{(t)} + \frac{\eta q}{nm} \cdot (2 \cdot \sqrt{\log(8mn/\delta)} \cdot \sigma_0 \sigma_p \sqrt{d/(n(p+s))} + 4\Psi^{(t)})^{q-1} \cdot 2\sigma_p^2 d \\
&\leq \Psi^{(t)} + \frac{\eta q}{nm} \cdot (4 \cdot \sqrt{\log(8mn/\delta)} \cdot \sigma_0 \sigma_p \sqrt{d/(n(p+s))})^{q-1} \cdot 2\sigma_p^2 d,
\end{aligned}$$

674 where the second inequality is due to the constraint $|\ell_i^{(t)}| \leq 1$, the third inequality is derived from
675 Lemmas G.1 and G.7, the fourth inequality is a consequence of the condition $d \geq 16Dn^2 \log(4n^2/\delta)$,
676 and the final inequality is a result of the inductive assumption (29). Summing over the sequence
677 $t = 0, 1, \dots, \tilde{T} - 1$, we obtain

$$\begin{aligned}
\Psi^{(\tilde{T})} &\leq \tilde{T} \cdot \frac{\eta q}{nm} \cdot (4 \cdot \sqrt{\log(8mn/\delta)} \cdot \sigma_0 \sigma_p \sqrt{d/(n(p+s))})^{q-1} \cdot 2\sigma_p^2 d \\
&\leq T_1^+ \cdot \frac{\eta q}{nm} \cdot (4 \cdot \sqrt{\log(8mn/\delta)} \cdot \sigma_0 \sigma_p \sqrt{d/(n(p+s))})^{q-1} \cdot 2\sigma_p^2 d \\
&\leq \frac{\sigma_0 \sigma_p \sqrt{d/(n(p+s))}}{2},
\end{aligned}$$

678 where the second inequality is justified by $\tilde{T} \leq T_1^+$ in our inductive argument. Hence, by induction,
679 we conclude that $\Psi^{(t)} \leq \sigma_0 \sigma_p \sqrt{d/(n(p+s))}/2$ for all $t \leq T_1^+$.

680 Next, we can assume, without loss of generality, that $j = 1$. Let $T_{1,1}$ represent the final time for t
681 within the interval $[0, T_1^+]$ such that $\max_r \gamma_{1,r}^{(t)} \leq 2$, given $\sigma_0 \leq \sqrt{n(p+s)/d}/\sigma_p$. For $t \leq T_{1,1}$,
682 we have $\max_{j,r,i} \{\rho_{j,r,i}^{(t)}\} = O(\sigma_0 \sigma_p \sqrt{d/(n(p+s))}) = O(1)$ and $\max_r \gamma_{1,r}^{(t)} \leq 2$. By applying
683 Lemmas H.8 and H.10, we deduce that $F_{-1}(\mathbf{W}_{-1}^{(t)}, \tilde{\mathbf{x}}_i), F_{+1}(\mathbf{W}_{+1}^{(t)}, \tilde{\mathbf{x}}_i) = O(1)$ for all i with $y_i = 1$.
684 Consequently, there exists a positive constant C_1 such that $-\ell_i^{(t)} \geq C_1$ for all i with $y_i = 1$.

685 By (10), for $t \leq T_{1,1}$ we have

$$\begin{aligned}
\gamma_{1,r}^{(t+1)} &= \gamma_{1,r}^{(t)} - \frac{\eta}{nm} \cdot \sum_{i=1}^n \ell_i^{(t)} \cdot \sigma'(\tilde{y}_i \cdot \langle \mathbf{w}_{1,r}^{(0)}, \boldsymbol{\mu} \rangle + \tilde{y}_i \cdot \gamma_{1,r}^{(t)}) \cdot \tilde{y}_i \|\boldsymbol{\mu}\|_2^2 \\
&\geq \gamma_{1,r}^{(t)} + \frac{C_1 \eta}{nm} \cdot \sum_{y_i=1} \sigma'(y_i \Xi \cdot \langle \mathbf{w}_{1,r}^{(0)}, \boldsymbol{\mu} \rangle + y_i \Xi \cdot \gamma_{1,r}^{(t)}) \cdot \frac{p-s}{p+s} \|\boldsymbol{\mu}\|_2^2.
\end{aligned}$$

686 Denote $\hat{\gamma}_{1,r}^{(t)} = \gamma_{1,r}^{(t)} + \langle \mathbf{w}_{1,r}^{(0)}, \boldsymbol{\mu} \rangle$ and let $A^{(t)} = \max_r \hat{\gamma}_{1,r}^{(t)}$. Then we have

$$\begin{aligned}
A^{(t+1)} &\geq A^{(t)} + \frac{C_1 \eta}{nm} \cdot \sum_{y_i=1} \sigma'(\Xi A^{(t)}) \cdot \Xi \|\boldsymbol{\mu}\|_2^2 \\
&\geq A^{(t)} + \frac{C_1 \eta q \|\boldsymbol{\mu}\|_2^2}{4m} \left[\Xi A^{(t)} \right]^{q-1} \Xi \\
&\geq \left(1 + \frac{C_1 \eta q \|\boldsymbol{\mu}\|_2^2}{4m} [A^{(0)}]^{q-2} \Xi^q \right) A^{(t)} \\
&\geq \left(1 + \frac{C_1 \eta q \sigma_0^{q-2} \|\boldsymbol{\mu}\|_2^q \Xi^q}{2^q m} \right) A^{(t)},
\end{aligned}$$

687 where the second inequality arises from the lower bound on the quantity of positive data as established
688 in Lemma G.4, the third inequality is a result of the increasing nature of the sequence $A^{(t)}$, and
689 the final inequality is derived from $A^{(0)} = \max_r \langle \mathbf{w}_{1,r}^{(0)}, \boldsymbol{\mu} \rangle \geq \sigma_0 \|\boldsymbol{\mu}\|_2 / 2$, as proven in Lemma G.7.
690 Consequently, the sequence $A^{(t)}$ exhibits exponential growth, and we can express it as

$$\begin{aligned}
A^{(t)} &\geq \left(1 + \frac{C_1 \eta q \sigma_0^{q-2} \|\boldsymbol{\mu}\|_2^q \Xi^q}{2^q m} \right)^t A^{(0)} \\
&\geq \exp \left(\frac{C_1 \eta q \sigma_0^{q-2} \|\boldsymbol{\mu}\|_2^q \Xi^q}{2^{q+1} m} t \right) A^{(0)} \\
&\geq \exp \left(\frac{C_1 \eta q \sigma_0^{q-2} \|\boldsymbol{\mu}\|_2^q \Xi^q}{2^{q+1} m} t \right) \frac{\sigma_0 \|\boldsymbol{\mu}\|_2}{2},
\end{aligned}$$

691 where the second inequality is justified by the relation $1 + z \geq \exp(z/2)$ for $z \leq 2$ and our specific
692 conditions on η and σ_0 as listed in Condition 3.1. The last inequality is a consequence of Lemma G.7
693 and the definition of $A^{(0)}$. Thus, $A^{(t)}$ will attain the value of 2 within T_1 iterations, defined as

$$T_1 = \frac{\log(6/\sigma_0 \|\boldsymbol{\mu}\|_2) 2^{q+1} m}{C_1 \eta q \sigma_0^{q-2} \|\boldsymbol{\mu}\|_2^q \Xi^q}.$$

694 Since $\max_r \gamma_{1,r}^{(t)} \geq A^{(t)} - 1$, $\max_r \gamma_{1,r}^{(t)}$ will reach 2 within T_1 iterations. Next, we can confirm that

$$T_1 \leq \frac{nm\eta^{-1}\sigma_0^{2-q}\sigma_p^{-q}d^{-q/2}(n(p+s))^{(q-2)/2}}{2^{q+5}q[4\log(8mn/\delta)]^{(q-1)/2}} = T_1^+/2,$$

695 where the inequality is consistent with our SNR condition in (27). Therefore, by the definition of
696 $T_{1,1}$, we deduce that $T_{1,1} \leq T_1 \leq T_1^+/2$, utilizing the non-decreasing property of γ . The proof
697 for $j = -1$ follows a similar logic, leading us to the conclusion that $\max_r \gamma_{-1,r}^{(T_1-1)} \geq 2$ while
698 $T_{1,-1} \leq T_1 \leq T_1^+/2$, thereby completing the proof.

699 □

700 I.2 Second stage: convergence analysis

701 After the first stage and at time step T_1 we know that:

$$\mathbf{w}_{j,r}^{(T_1)} = \mathbf{w}_{j,r}^{(0)} + j \cdot \gamma_{j,r}^{(T_1)} \cdot \frac{\boldsymbol{\mu}}{\|\boldsymbol{\mu}\|_2^2} + \sum_{i=1}^n \bar{\rho}_{j,r,i}^{(T_1)} \cdot \frac{\boldsymbol{\xi}_i}{\|\boldsymbol{\xi}_i\|_2^2} + \sum_{i=1}^n \underline{\rho}_{j,r,i}^{(T_1)} \cdot \frac{\boldsymbol{\xi}_i}{\|\boldsymbol{\xi}_i\|_2^2}.$$

702 And at the beginning of the second stage, we have following property holds:

- 703 • $\max_r \gamma_{j,r}^{(T_1)} \geq 2, \forall j \in \{\pm 1\}$.
- 704 • $\max_{j,r,i} |\rho_{j,r,i}^{(T_1)}| \leq \hat{\beta}$ where $\hat{\beta} = \sigma_0 \sigma_p \sqrt{d/(n(p+s))}/2$.

705 Lemma E.1 implies that the learned feature $\gamma_{j,r}^{(t)}$ will not get worse, i.e., for $t \geq T_1$, we have that
 706 $\gamma_{j,r}^{(t+1)} \geq \gamma_{j,r}^{(t)}$, and therefore $\max_r \gamma_{j,r}^{(t)} \geq 2$. Now we choose \mathbf{W}^* as follows:

$$\mathbf{w}_{j,r}^* = \mathbf{w}_{j,r}^{(0)} + 2qm \log(2q/\epsilon) \cdot j \cdot \frac{\boldsymbol{\mu}}{\|\boldsymbol{\mu}\|_2^2}.$$

707 While the context of CNN presents subtle differences from the scenario described in CNN [15], we
 708 can adapt the same analytical approach to derive the following two lemmas:

709 **Lemma I.3** ([15]). *Under the same conditions as Theorem 3.2, we have that $\|\mathbf{W}^{(T_1)} - \mathbf{W}^*\|_F \leq$
 710 $\tilde{O}(m^{3/2}\|\boldsymbol{\mu}\|_2^{-1})$.*

711 **Lemma I.4** ([15]). *Under the same conditions as Theorem 3.2, we have that*

$$\|\mathbf{W}^{(t)} - \mathbf{W}^*\|_F^2 - \|\mathbf{W}^{(t+1)} - \mathbf{W}^*\|_F^2 \geq (2q-1)\eta L_S(\mathbf{W}^{(t)}) - \eta\epsilon$$

712 for all $T_1 \leq t \leq T^*$.

713 **Lemma I.5** (Restatement of Lemma E.3). *Under the same conditions as Theorem 3.2, let $T =$
 714 $T_1 + \left\lfloor \frac{\|\mathbf{W}^{(T_1)} - \mathbf{W}^*\|_F^2}{2\eta\epsilon} \right\rfloor = T_1 + \tilde{O}(m^3\eta^{-1}\epsilon^{-1}\|\boldsymbol{\mu}\|_2^{-2})$. Then we have $\max_{j,r,i} |\rho_{j,r,i}^{(t)}| \leq 2\hat{\beta} =$
 715 $\sigma_0\sigma_p\sqrt{d/(n(p+s))}$ for all $T_1 \leq t \leq T$. Besides,*

$$\frac{1}{t - T_1 + 1} \sum_{s=T_1}^t L_S(\mathbf{W}^{(s)}) \leq \frac{\|\mathbf{W}^{(T_1)} - \mathbf{W}^*\|_F^2}{(2q-1)\eta(t - T_1 + 1)} + \frac{\epsilon}{2q-1}$$

716 for all $T_1 \leq t \leq T$, and we can find an iterate with training loss smaller than ϵ within T iterations.

717 *Proof of Lemma I.5.* We adapt the convergence proof for CNN[15] to extend the analysis to GNN.
 718 By invoking Lemma I.4, for any given time interval $t \in [T_1, T]$, we can deduce that

$$\|\mathbf{W}^{(s)} - \mathbf{W}^*\|_F^2 - \|\mathbf{W}^{(s+1)} - \mathbf{W}^*\|_F^2 \geq (2q-1)\eta L_S(\mathbf{W}^{(s)}) - \eta\epsilon,$$

719 which is valid for $s \leq t$. Summing over this interval, we arrive at

$$\sum_{s=T_1}^t L_S(\mathbf{W}^{(s)}) \leq \frac{\|\mathbf{W}^{(T_1)} - \mathbf{W}^*\|_F^2 + \eta\epsilon(t - T_1 + 1)}{(2q-1)\eta}. \quad (30)$$

720 This inequality holds for all $T_1 \leq t \leq T$. Dividing both sides of (30) by $(t - T_1 + 1)$, we obtain

$$\frac{1}{t - T_1 + 1} \sum_{s=T_1}^t L_S(\mathbf{W}^{(s)}) \leq \frac{\|\mathbf{W}^{(T_1)} - \mathbf{W}^*\|_F^2}{(2q-1)\eta(t - T_1 + 1)} + \frac{\epsilon}{2q-1}.$$

721 By setting $t = T$, we find that

$$\frac{1}{T - T_1 + 1} \sum_{s=T_1}^T L_S(\mathbf{W}^{(s)}) \leq \frac{\|\mathbf{W}^{(T_1)} - \mathbf{W}^*\|_F^2}{(2q-1)\eta(T - T_1 + 1)} + \frac{\epsilon}{2q-1} \leq \frac{3\epsilon}{2q-1} < \epsilon,$$

722 where we utilize the condition that $q > 2$ and the specific choice of $T = T_1 + \left\lfloor \frac{\|\mathbf{W}^{(T_1)} - \mathbf{W}^*\|_F^2}{2\eta\epsilon} \right\rfloor$.
 723 Since the mean value is less than ϵ , it follows that there must exist a time interval $T_1 \leq t \leq T$ for
 724 which $L_S(\mathbf{W}^{(t)}) < \epsilon$.

725 Finally, we aim to demonstrate that $\max_{j,r,i} |\rho_{j,r,i}^{(t)}| \leq 2\hat{\beta}$ holds for all $t \in [T_1, T]$. By inserting
 726 $T = T_1 + \left\lfloor \frac{\|\mathbf{W}^{(T_1)} - \mathbf{W}^*\|_F^2}{2\eta\epsilon} \right\rfloor$ into equation (30), we obtain

$$\sum_{s=T_1}^T L_S(\mathbf{W}^{(s)}) \leq \frac{2\|\mathbf{W}^{(T_1)} - \mathbf{W}^*\|_F^2}{(2q-1)\eta} = \tilde{O}(\eta^{-1}m^3\|\boldsymbol{\mu}\|_2^2), \quad (31)$$

727 where the inequality is a consequence of $\|\mathbf{W}^{(T_1)} - \mathbf{W}^*\|_F \leq \tilde{O}(m^{3/2}\|\boldsymbol{\mu}\|_2^{-1})$ as shown in Lemma I.3.

728 Let's define $\Psi^{(t)} = \max_{j,r,i} |\rho_{j,r,i}^{(t)}|$. We will employ induction to prove $\Psi^{(t)} \leq 2\hat{\beta}$ for all $t \in [T_1, T]$.

729 At $t = T_1$, by the definition of $\hat{\beta}$, it is clear that $\Psi^{(T_1)} \leq \hat{\beta} \leq 2\hat{\beta}$.

730 Assuming that there exists $\tilde{T} \in [T_1, T]$ such that $\Psi^{(t)} \leq 2\hat{\beta}$ for all $t \in [T_1, \tilde{T} - 1]$, we can consider

731 $t \in [T_1, \tilde{T} - 1]$. Using the expression:

$$\begin{aligned} \rho_{j,r,i}^{(t+1)} &= \rho_{j,r,i}^{(t)} - \frac{\eta}{nm} \cdot \sum_{k \in \mathcal{N}(i)} D_k^{-1} \ell_k^{(t)} \\ &\sigma' \left(\langle \mathbf{w}_{j,r}^{(0)}, \tilde{\xi}_k \rangle + \sum_{i'=1}^n \tilde{\rho}_{j,r,i'}^{(t)} \frac{\langle \xi_{i'}, \tilde{\xi}_k \rangle}{\|\xi_{i'}\|_2^2} + \sum_{i'=1}^n \rho_{j,r,i'}^{(t)} \frac{\langle \xi_{i'}, \tilde{\xi}_k \rangle}{\|\xi_{i'}\|_2^2} \right) \cdot \|\xi_i\|_2^2 \end{aligned} \quad (32)$$

732 we can proceed to analyze:

$$\begin{aligned} \Psi^{(t+1)} &\leq \Psi^{(t)} + \max_{j,r,i} \left\{ \frac{\eta}{nm} \cdot \sum_{k \in \mathcal{N}(i)} D_k^{-1} |\ell_k^{(t)}| \cdot \sigma' \left(\langle \mathbf{w}_{j,r}^{(0)}, \tilde{\xi}_k \rangle + 2 \sum_{i'=1}^n \Psi^{(t)} \cdot \frac{|\langle \xi_{i'}, \tilde{\xi}_k \rangle|}{\|\xi_{i'}\|_2^2} \right) \cdot \|\xi_i\|_2^2 \right\} \\ &= \Psi^{(t)} + \max_{j,r,i} \left\{ \frac{\eta}{nm} \cdot \sum_{k \in \mathcal{N}(i)} D_k^{-1} |\ell_k^{(t)}| \cdot \right. \\ &\quad \left. \sigma' \left(\langle \mathbf{w}_{j,r}^{(0)}, \tilde{\xi}_k \rangle + 2 \sum_{i' \neq k'}^n \Psi^{(t)} \cdot D_k^{-1} \sum_{k' \in \mathcal{N}(k)} \frac{|\langle \xi_{i'}, \xi_{k'} \rangle|}{\|\xi_{i'}\|_2^2} \right) \cdot \|\xi_i\|_2^2 \right\} \\ &\leq \Psi^{(t)} + \frac{\eta q}{nm} \cdot \max_i \sum_{k \in \mathcal{N}(i)} D_k^{-1} |\ell_k^{(t)}| \cdot \left[2 \cdot \sqrt{\log(8mn/\delta)} \cdot \sigma_0 \sigma_p \sqrt{d/(n(p+s))} \right. \\ &\quad \left. + \left(2 + \frac{4n\sigma_p^2 \cdot \sqrt{d \log(4n^2/\delta)}}{\sigma_p^2 d/2} \right) \cdot \Psi^{(t)} \right]^{q-1} \cdot 2\sigma_p^2 d \\ &\leq \Psi^{(t)} + \frac{\eta q}{nm} \cdot \max_i \sum_{k \in \mathcal{N}(i)} D_k^{-1} |\ell_k^{(t)}| \cdot \\ &\quad \left(2 \cdot \sqrt{\log(8mn/\delta)} \cdot \sigma_0 \sigma_p \sqrt{d/(n(p+s))} + 4 \cdot \Psi^{(t)} \right)^{q-1} \cdot 2\sigma_p^2 d. \end{aligned}$$

733 The second inequality is derived from Lemmas G.1 and G.7, while the final inequality is based on the
734 assumption that $d \geq 16n^2 \log(4n^2/\delta)$. By taking a telescoping sum, we can express the following:

$$\begin{aligned} \Psi^{(T)} &\stackrel{(i)}{\leq} \Psi^{(T_1)} + \frac{\eta q}{nm} \sum_{s=T_1}^{\tilde{T}-1} \max_i \sum_{k \in \mathcal{N}(i)} D_k^{-1} |\ell_k^{(s)}| \tilde{O}(\sigma_p^2 d) \hat{\beta}^{q-1} \\ &\stackrel{(ii)}{\leq} \Psi^{(T_1)} + \frac{\eta q}{nm} \tilde{O}(\sigma_p^2 d) \hat{\beta}^{q-1} \sum_{s=T_1}^{\tilde{T}-1} \max_i \sum_{k \in \mathcal{N}(i)} D_k^{-1} \ell_k^{(s)} \\ &\stackrel{(iii)}{\leq} \Psi^{(T_1)} + \tilde{O}(\eta m^{-1} \sigma_p^2 d) \hat{\beta}^{q-1} \sum_{s=T_1}^{\tilde{T}-1} L_S(\mathbf{W}^{(s)}) \\ &\stackrel{(iv)}{\leq} \Psi^{(T_1)} + \tilde{O}(m^2 \text{SNR}^{-2}) \hat{\beta}^{q-1} \\ &\stackrel{(v)}{\leq} \hat{\beta} + \tilde{O}(m^2 n^{2/q} (n(p+s))^{1-2/q} \hat{\beta}^{q-2}) \hat{\beta} \\ &\stackrel{(vi)}{\leq} 2\hat{\beta}, \end{aligned}$$

735 where (i) follows from our induction assumption that $\Psi^{(t)} \leq 2\hat{\beta}$, (ii) is derived from the relationship

736 $|\ell'| \leq \ell$, (iii) is obtained by the sum of $\max_i \sum_{k \in \mathcal{N}(i)} D_k^{-1} \leq \sum_i \ell_i^{(s)} = nL_S(\mathbf{W}^{(s)})$, (iv) is

737 due to the summation of $\sum_{s=T_1}^{\tilde{T}-1} L_S(\mathbf{W}^{(s)}) \leq \sum_{s=T_1}^T L_S(\mathbf{W}^{(s)}) = \tilde{O}(\eta^{-1} m^3 \|\boldsymbol{\mu}\|_2^2)$ as shown in

738 (31), (v) is based on the condition $n\text{SNR}^q \cdot (n(p+s))^{q/2-1} \geq \tilde{\Omega}(1)$, and (vi) follows from the

739 definition of $\hat{\beta} = \sigma_0 \sigma_p \sqrt{d/(n(p+s))}/2$ and $\tilde{O}(m^2 n^{2/q} (n(p+s))^{1-2/q} \hat{\beta}^{q-2}) = \tilde{O}(m^2 n^{2/q} (n(p+s))^{1-2/q} (\sigma_0 \sigma_p \sqrt{d/(n(p+s))})^{q-2}) \leq 1$.

741 Thus, we conclude that $\Psi(\tilde{T}) \leq 2\hat{\beta}$, completing the induction and establishing the desired result. \square

742 I.3 Population loss

743 Consider a new data point (\mathbf{x}, y) drawn from the SNM-SBM distribution. Without loss of generality,
744 we suppose that the first patch is the signal patch and the second patch is the noise patch, i.e.,
745 $\mathbf{x} = [y \cdot \boldsymbol{\mu}, \boldsymbol{\xi}]$. Moreover, by the signal-noise decomposition, the learned neural network has
746 parameter:

$$\mathbf{w}_{j,r}^{(t)} = \mathbf{w}_{j,r}^{(0)} + j \cdot \gamma_{j,r}^{(t)} \cdot \frac{\boldsymbol{\mu}}{\|\boldsymbol{\mu}\|_2} + \sum_{i=1}^n \bar{\rho}_{j,r,i}^{(t)} \cdot \frac{\boldsymbol{\xi}_i}{\|\boldsymbol{\xi}_i\|_2} + \sum_{i=1}^n \underline{\rho}_{j,r,i}^{(t)} \cdot \frac{\boldsymbol{\xi}_i}{\|\boldsymbol{\xi}_i\|_2}$$

747 for $j \in \{\pm 1\}$ and $r \in [m]$.

748 Although the framework of CNN diverges in certain nuances from the situation of CNN outlined in
749 [15], we are able to employ a similar analytical methodology to deduce the subsequent two lemmas:

750 **Lemma I.6.** *Under the same conditions as Theorem 3.2, we have that $\max_{j,r} |\langle \mathbf{w}_{j,r}^{(t)}, \tilde{\boldsymbol{\xi}}_i \rangle| \leq 1/2$ for
751 all $0 \leq t \leq T$, and $i \in [n]$.*

752 **Lemma I.7.** *Under the same conditions as Theorem 3.2, with probability at least $1 - 4mT \cdot$
753 $\exp(-C_2^{-1} \sigma_0^{-2} \sigma_p^{-2} d^{-1} n(p+s))$, we have that $\max_{j,r} |\langle \mathbf{w}_{j,r}^{(t)}, \tilde{\boldsymbol{\xi}} \rangle| \leq 1/2$ for all $0 \leq t \leq T$, where
754 $C_2 = \tilde{O}(1)$.*

755 **Lemma I.8** (Restatement of Lemma E.4). *Let T be defined in Lemma E.2 respectively. Under
756 the same conditions as Theorem 3.2, for any $0 \leq t \leq T$ with $L_S(\mathbf{W}^{(t)}) \leq 1$, it holds that
757 $L_{\mathcal{D}}(\mathbf{W}^{(t)}) \leq c_1 \cdot L_S(\mathbf{W}^{(t)}) + \exp(-c_2 n^2)$.*

758 *Proof of Lemma I.8.* Consider the occurrence of event \mathcal{E} , defined as the condition under which
759 Lemma I.7 is satisfied. We can then express the loss $L_{\mathcal{D}}(\mathbf{W}^{(t)})$ as a sum of two components:

$$\mathbb{E}[\ell(yf(\mathbf{W}^{(t)}, \tilde{\mathbf{x}}))] = \underbrace{\mathbb{E}[\mathbb{1}(\mathcal{E})\ell(yf(\mathbf{W}^{(t)}, \tilde{\mathbf{x}}))]}_{\text{Term } I_1} + \underbrace{\mathbb{E}[\mathbb{1}(\mathcal{E}^c)\ell(yf(\mathbf{W}^{(t)}, \tilde{\mathbf{x}}))]}_{\text{Term } I_2}. \quad (33)$$

760 Next, we proceed to establish bounds for I_1 and I_2 .

761 **Bounding I_1 :** Given that $L_S(\mathbf{W}^{(t)}) \leq 1$, there must be an instance $(\tilde{\mathbf{x}}_i, y_i)$ for which
762 $\ell(y_i f(\mathbf{W}^{(t)}, \tilde{\mathbf{x}}_i)) \leq L_S(\mathbf{W}^{(t)}) \leq 1$, leading to $y_i f(\mathbf{W}^{(t)}, \tilde{\mathbf{x}}_i) \geq 0$. Hence, we obtain:

$$\exp(-y_i f(\mathbf{W}^{(t)}, \tilde{\mathbf{x}}_i)) \stackrel{(i)}{\leq} 2 \log(1 + \exp(-y_i f(\mathbf{W}^{(t)}, \tilde{\mathbf{x}}_i))) = 2\ell(y_i f(\mathbf{W}^{(t)}, \tilde{\mathbf{x}}_i)) \leq 2L_S(\mathbf{W}^{(t)}), \quad (34)$$

763 where (i) follows from the inequality $z \leq 2 \log(1+z), \forall z \leq 1$. If event \mathcal{E} occurs, we deduce:

$$\begin{aligned} |yf(\mathbf{W}^{(t)}, \tilde{\mathbf{x}}^{(2)}) - y_i f(\mathbf{W}^{(t)}, \tilde{\mathbf{x}}_i^{(2)})| &\leq \frac{1}{m} \sum_{j,r} \sigma(\langle \mathbf{w}_{j,r}, \tilde{\boldsymbol{\xi}}_i \rangle) + \frac{1}{m} \sum_{j,r} \sigma(\langle \mathbf{w}_{j,r}, \tilde{\boldsymbol{\xi}} \rangle) \\ &\leq 1. \end{aligned} \quad (35)$$

764 Here, $f(\mathbf{W}^{(t)}, \tilde{\mathbf{x}}^{(2)})$ refers to the input $\tilde{\mathbf{x}} = [0, \tilde{\mathbf{x}}^{(2)}]$. The second inequality is justified by Lemmas I.7
765 and I.6. Consequently, we have:

$$\begin{aligned}
I_1 &\leq \mathbb{E}[\mathbb{1}(\mathcal{E}) \exp(-yf(\mathbf{W}^{(t)}, \tilde{\mathbf{x}}))] \\
&= \mathbb{E}[\mathbb{1}(\mathcal{E}) \exp(-y_i f(\mathbf{W}^{(t)}, \tilde{\mathbf{x}}^{(1)})) \exp(-y_i f(\mathbf{W}^{(t)}, \tilde{\mathbf{x}}^{(2)}))] \\
&\leq 2e \cdot C \cdot \mathbb{E}[\mathbb{1}(\mathcal{E}) \exp(-y_i f(\mathbf{W}^{(t)}, \tilde{\mathbf{x}}_i^{(1)})) \exp(-y_i f(\mathbf{W}^{(t)}, \tilde{\mathbf{x}}_i^{(2)}))] \\
&\leq 2e \cdot \mathbb{E}[\mathbb{1}(\mathcal{E}) L_{\mathcal{S}}(\mathbf{W}^{(t)})],
\end{aligned}$$

766 where the inequalities follow from the properties of cross-entropy loss, (35), Lemma G.4, and (34).
767 The constant c_1 encapsulates the factors in the derivation.

768 **Estimating I_2 :** We now turn our attention to the second term I_2 . By selecting an arbitrary training
769 data point $(\mathbf{x}_{i'}, y_{i'})$ with $y_{i'} = y$, we can derive the following:

$$\begin{aligned}
\ell(yf(\mathbf{W}^{(t)}, \tilde{\mathbf{x}})) &\leq \log(1 + \exp(F_{-y}(\mathbf{W}^{(t)}, \tilde{\mathbf{x}}))) \\
&\leq 1 + F_{-y}(\mathbf{W}^{(t)}, \tilde{\mathbf{x}}) \\
&= 1 + \frac{1}{m} \sum_{j=-y, r \in [m]} \sigma(\langle \mathbf{w}_{j,r}^{(t)}, \tilde{y}\boldsymbol{\mu} \rangle) + \frac{1}{m} \sum_{j=-y, r \in [m]} \sigma(\langle \mathbf{w}_{j,r}^{(t)}, \tilde{\boldsymbol{\xi}} \rangle) \\
&\leq 1 + F_{-y_i}(\mathbf{W}_{-y_{i'}}, \tilde{\mathbf{x}}_{i'}) + \frac{1}{m} \sum_{j=-y, r \in [m]} \sigma(\langle \mathbf{w}_{j,r}^{(t)}, \tilde{\boldsymbol{\xi}} \rangle) \\
&\leq 2 + \frac{1}{m} \sum_{j=-y, r \in [m]} \sigma(\langle \mathbf{w}_{j,r}^{(t)}, \tilde{\boldsymbol{\xi}} \rangle) \\
&\leq 2 + \tilde{O}((\sigma_0 \sqrt{d})^q) \|\tilde{\boldsymbol{\xi}}\|^q, \tag{36}
\end{aligned}$$

770 where the inequalities follow from the properties of the cross-entropy loss and the constraints defined
771 in Lemma H.8. The last inequality is a result of the boundedness of the inner product with $\tilde{\boldsymbol{\xi}}$.
772 Continuing, we have:

$$\begin{aligned}
I_2 &\leq \sqrt{\mathbb{E}[\mathbb{1}(\mathcal{E}^c)]} \cdot \sqrt{\mathbb{E}[\ell(yf(\mathbf{W}^{(t)}, \tilde{\mathbf{x}}))^2]} \\
&\leq \sqrt{\mathbb{P}(\mathcal{E}^c)} \cdot \sqrt{4 + \tilde{O}((\sigma_0 \sqrt{d})^{2q}) \mathbb{E}[\|\tilde{\boldsymbol{\xi}}\|_2^{2q}]} \\
&\leq \exp \left[-\tilde{\Omega} \left(\frac{\sigma_0^{-2} \sigma_p^{-2}}{d^{-1} n(p+s)} \right) + \text{polylog}(d) \right] \\
&\leq \exp(-c_1 n^2),
\end{aligned}$$

773 where c_1 is a constant, the first inequality is by Cauchy-Schwartz inequality, the second inequality is
774 by (36), the third inequality is by Lemma I.7 and the fact that $\sqrt{4 + \tilde{O}((\sigma_0 \sqrt{d})^{2q}) \mathbb{E}[\|\tilde{\boldsymbol{\xi}}\|_2^{2q}]} =$
775 $O(\text{poly}(d))$, and the last inequality is by our condition $\sigma_0 \leq \tilde{O}(m^{-2/(q-2)} n^{-1}) \cdot$
776 $(\sigma_p \sqrt{d}/(n(p+s)))^{-1}$ in Condition 3.1. Plugging the bounds of I_1, I_2 completes the proof. \square

777 J Additional Experimental Procedures and Results

778 J.1 Dataset in Node Classification

779 In Figure 1, we execute node classification experiments on three frequently used citation networks:
780 Cora, Citeseer, and Pubmed [1]. Detailed information about these datasets is provided below and
781 summarized in Table 1.

- 782 • The Cora dataset includes 2,708 scientific publications, each categorized into one of seven
783 classes, connected by 5,429 links. Each publication is represented by a binary word vector,

Table 1: Details of Datasets

Dataset	Nodes	Edges	Classes	Features	Train/Val/Test
Cora	2,708	5,429	7	1,433	0.05/0.18/0.37
Citeseer	3,327	4,732	6	3,703	0.04/0.15/0.30
Pubmed	19,717	44,338	3	500	0.003/0.03/0.05

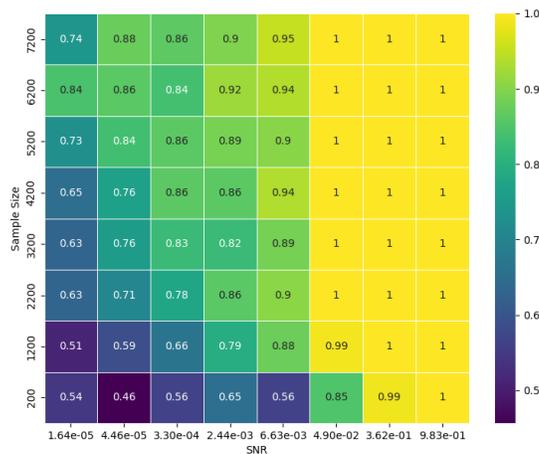


Figure 7: Test accuracy heatmap for GCNs after training.

784 which denotes the presence or absence of a corresponding word from a dictionary of 1,433
785 unique words.

- 786 • The Citeseer dataset comprises 3,312 scientific publications, each classified into one of six
787 classes, connected by 4,732 links. Each publication is represented by a binary word vector,
788 indicating the presence or absence of a corresponding word from a dictionary that includes
789 3,703 unique words.
- 790 • The Pubmed Diabetes dataset includes 19,717 scientific publications related to diabetes,
791 drawn from the PubMed database and classified into one of three classes. The citation
792 network is made up of 44,338 links. Each publication is represented by a TF-IDF weighted
793 word vector from a dictionary consisting of 500 unique words.

794 J.2 Phase transition in GCN

795 In Figure 5, we illustrated the variance in test accuracy between CNN and GCN within a chosen range
796 of SNR and sample numbers, where GCN was shown to achieve near-perfect test accuracy. Here,
797 we broaden the SNR range towards the smaller end and display the corresponding phase diagram
798 of GCN in Figure 7. When the SNR is exceedingly small, we observe that GCNs return lower test
799 accuracy, suggesting the possibility of a phase transition in the test accuracy of GCNs.

800 J.3 Software and hardware

801 We implement our methods with PyTorch. For the software and hardware configurations, we ensure
802 the consistent environments for each datasets. We run all the experiments on Linux servers with
803 NVIDIA V100 graphics cards with CUDA 11.2.