

---

# HERTA: A High-Efficiency and Rigorous Training Algorithm for Unfolded Graph Neural Networks

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 As a variant of Graph Neural Networks (GNNs), Unfolded GNNs offer enhanced  
2 interpretability and flexibility over traditional designs. Nevertheless, they still suffer  
3 from scalability challenges when it comes to the training cost. Although many  
4 methods have been proposed to address the scalability issues, they mostly focus  
5 on per-iteration efficiency, without worst-case convergence guarantees. Moreover,  
6 those methods typically add components to or modify the original model, thus  
7 possibly breaking the interpretability of Unfolded GNNs. In this paper, we propose  
8 HERTA: a High-Efficiency and Rigorous Training Algorithm for Unfolded GNNs  
9 that accelerates the whole training process, achieving a nearly-linear time worst-  
10 case training guarantee. Crucially, HERTA converges to the optimum of the original  
11 model, thus preserving the interpretability of Unfolded GNNs. Additionally, as a  
12 byproduct of HERTA, we propose a new spectral sparsification method applicable  
13 to normalized and regularized graph Laplacians that ensures tighter bounds for our  
14 algorithm than existing spectral sparsifiers do. Experiments on real-world datasets  
15 verify the superiority of HERTA as well as its adaptability to various loss functions  
16 and optimizers.

## 17 1 Introduction

18 Graph Neural Networks (GNNs) have become a powerful modern tool for handling graph data  
19 because of their strong representational capabilities and ability to explore the relationships between  
20 data points [22, 2, 42]. Like many deep learning models, GNNs are generally designed by heuristics  
21 and experience, which makes analyzing and understanding them a difficult task. *Unfolded GNNs*  
22 [45] are a type of GNNs that are rigorously derived from an optimization problem, which makes their  
23 inductive bias and training dynamics easier to interpret.

24 Despite their enhanced interpretability, training unfolded models can be extremely expensive, espe-  
25 cially when the graph is large and (relatively) dense. This issue can come from two aspects:

- 26 1) *Slow iterations*: in every iteration the whole graph must be fed into the model, and the interactive  
27 nature of graphs prevents trivially utilizing techniques like mini-batch training;
- 28 2) *Slow convergence*: the training convergence rate is related to the connectivity of the graph and  
29 the well-conditionedness of the node features, and the model converges slowly when the data is  
30 ill-conditioned.

31 Many methods have been proposed to address the high cost of training unfolded models [5, 7, 47,  
32 14, 16, 20], primarily by using graph sampling to enable mini-batch training schemes that reduce the  
33 per-iteration cost (Issue 1). However, these approaches typically require distorting the underlying  
34 optimization objective explicitly or implicitly, thus diminishing the rigorous and interpretable nature  
35 of Unfolded GNNs. Moreover, existing works have mostly focused on per-iteration efficiency, while

36 the convergence rate of the optimization in training Unfolded GNNs (Issue 2) remains un-addressed,  
37 leading to methods that are not robust to ill-conditioned data.

38 We propose HERTA: a High-Efficiency and Rigorous Training Algorithm for Unfolded GNNs,  
39 which is an algorithmic framework designed to address both Issues 1 and 2, while at the same  
40 time preserving the rigorous and interpretable nature of Unfolded GNNs. Unlike many existing  
41 methods that require changing the GNN model or defining a new objective, HERTA converges to  
42 the optimum of the original training problem, and thus preserves the interpretability of the unfolded  
43 models. Moreover, HERTA uses a specialized preconditioner to ensure fast linear convergence to the  
44 optimum, requiring only a logarithmic number of passes over the data. We show empirically that  
45 HERTA can be used to accelerate training for a variety of GNN objectives, as an extension of popular  
46 optimization techniques.

47 A key novelty of our algorithm lies in the construction of a preconditioner for the GNN objective,  
48 which accelerates the convergence rate. To construct this preconditioner, we design a spectral  
49 sparsification method for approximating the squared inverse of a normalized and regularized Laplacian  
50 matrix, by relying on an extension of the notion of effective resistance to overcome the limitations of  
51 existing graph sparsification tools. Below, we present an informal statement of our main result.

52 **Theorem 1.1** (Informal version of Theorem 5.1). *HERTA solves the  $\lambda$ -regularized Unfolded GNN*  
53 *objective (4.2) with  $n$  nodes,  $m$  edges and  $d$ -dimensional node features to within accuracy  $\epsilon$  in time*  
54  *$\tilde{O}\left((m + nd) \left(\log \frac{1}{\epsilon}\right)^2 + d^3\right)$  as long as the number of large eigenvalues of the graph Laplacian is*  
55  *$O(n/\lambda^2)$ .*

56 In practice, the node feature dimensionality  $d$  is generally much smaller than the graph size  $m$ , in  
57 which case, the running time of HERTA is  $\tilde{O}\left((m + nd) \left(\log \frac{1}{\epsilon}\right)^2\right)$ . Notice that to describe a typical  
58 graph dataset with  $m$  edges,  $n$  nodes and  $d$ -dimensional node features, we need at least  $O(m + nd)$   
59 float or integer numbers. This shows the essential optimality of HERTA: its running time is of the  
60 same magnitude as reading the input data (up to logarithm factors).

61 The condition that the graph Laplacian has  $O(n/\lambda^2)$  large eigenvalues is not strictly necessary; it is  
62 used here to simplify the time complexity (see Theorem 5.1 for the full statement). This condition is  
63 also not particularly restrictive, since in most practical settings the GNN parameter  $\lambda$  is chosen as an  
64 absolute constant. At the same time, it is natural that the complexity of training a GNN depends on  
65 the number of large eigenvalues in the graph Laplacian, since this quantity can be interpreted as the  
66 effective dimension in the Laplacian energy function (4.1) used to define the Unfolded GNN (see  
67 more discussion in Section 5).

68 **Outline.** The paper is organized as follows. In Section 2 we introduce related work. In Section 3 we  
69 introduce the mathematical notation and concepts used in this paper. We define our problem setting in  
70 Section 4. In Section 5 we present and analyze our algorithm HERTA, and introduce the techniques  
71 that are used in proving the main result. We conduct experiments on real-world datasets and show the  
72 results in Section 6. Finally, we conclude with some potential future directions in Section 7.

## 73 2 Related Work

74 **Unfolded Graph Neural Networks.** Unlike conventional GNN models, which are designed  
75 mainly by heuristics, the forward layers of Unfolded GNNs are derived by explicitly optimizing a  
76 graph-regularized target function [48, 45, 26, 25, 46, 50, 1]. This optimization nature of Unfolded  
77 GNNs allows developing models in a more interpretable and controllable way. The presence of an  
78 explicit optimization problem allows designers to inject desired properties into the model and better  
79 understand it. This approach has been used to overcome known GNN issues such as over-smoothing  
80 [30, 24, 44] and sensitivity towards spurious edges [51, 49].

81 **Efficient Graph Neural Network Training.** In order to address the scalability issue of GNNs, many  
82 techniques have been proposed. There are two major directions in the exploration of the solutions  
83 of this issue. The first direction is to adopt sampling methods. These methods include node-wise  
84 sampling [18], layer-wise sampling [6] and subgraph-wise sampling [7, 20]. The second direction is  
85 to utilize embeddings from past iterations to enhance current iterations and to obtain more accurate  
86 representation with mini-batches and fewer forward layers [16, 14, 43].

87 As we indicated in Section 1, all of the above methods aim at addressing Issue 1 (*Slow Iterations*),  
 88 and none of them achieves a worst-case convergence guarantee without affecting the underlying  
 89 GNN model. In contrast, HERTA addresses both Issue 1 and also Issue 2 (*Slow Convergence*), and  
 90 possesses a theoretical guarantee on the training convergence rate while preserving the original target  
 91 model.

92 **Matrix Sketching and Subsampling.** Our techniques and analysis are closely related to matrix  
 93 sketching and subsampling, which are primarily studied in the area of Randomized Numerical Linear  
 94 Algebra (RandNLA, [41, 15, 27, 29]). Given a large matrix, by sketching or sampling we compress  
 95 it to a much smaller matrix with certain properties preserved, thus accelerating the solution by  
 96 conducting expensive computations on the smaller matrix. This type of methods lead to improved  
 97 randomized algorithms for tasks including low rank approximation [17, 10, 8], linear system solving  
 98 [32, 13], least squares regression [33, 28] and so on [9]. Our usage of these methods includes  
 99 constructing spectral sparsifiers of graphs, obtained by edge subsampling [35, 36, 23], which are  
 100 central in designing fast solvers for Laplacian linear systems [39, 21, 31].

### 101 3 Preliminaries

102 For a vector  $\mathbf{x}$ , we denote its  $\ell_2$  norm as  $\|\mathbf{x}\|$ . We use  $\text{diag}(\mathbf{x})$  to denote diagonalization of  $\mathbf{x}$  into a  
 103 matrix. For a matrix  $M$ , we use  $\|M\|_{\mathcal{F}}$  and  $\|M\|$  to denote its Frobenius norm and operator norm.  
 104 Let  $\sigma_{\max}(M)$  and  $\sigma_{\min}(M)$  be the largest and smallest singular value of  $M$  respectively, we denote  
 105 its condition number as  $\kappa(M) = \frac{\sigma_{\max}(M)}{\sigma_{\min}(M)}$ . We use  $\text{nnz}(M)$  to denote the number of non-zero  
 106 entries of  $M$ . For two positive semidefinite (PSD) matrices  $\Sigma$  and  $\tilde{\Sigma}$ , if there exists  $\epsilon \in (0, 1)$  such  
 107 that  $(1 - \epsilon)\tilde{\Sigma} \preceq \Sigma \preceq (1 + \epsilon)\tilde{\Sigma}$ , then we say  $\Sigma \approx_{\epsilon} \tilde{\Sigma}$ , where  $\preceq$  refers to Loewner order. We use  
 108  $\delta_u$  to denote a vector with  $u$ -th entry equal to 1 and all other entries equal to 0. We also use  $\delta_{u,v}$  to  
 109 represent the  $v$ -th entry of  $\delta_u$ , i.e.  $\delta_{u,v} = \begin{cases} 1 & u = v \\ 0 & u \neq v \end{cases}$ .

110 For a function  $\ell : \mathbb{R}^d \rightarrow \mathbb{R}$  that is bounded from bottom, a point  $\mathbf{w}_0$  is called a solution of  $\ell$  with  $\epsilon$   
 111 error rate if it satisfies  $\ell(\mathbf{w}_0) \leq (1 + \epsilon)\ell^*$  where  $\ell^* = \inf_{\mathbf{w}} \ell(\mathbf{w})$ .

112 We adopt big-O notations in the time complexity analysis. Moreover, we use the notation  $\tilde{O}(\cdot)$  to  
 113 hide polynomial logarithmic factors of the problem size. As an example,  $O(\log(n))$  can be expressed  
 114 as  $\tilde{O}(1)$ .

115 For a matrix  $\Pi \in \mathbb{R}^{s \times n}$ , we call it a Gaussian sketch if its entries are i.i.d.  $\mathcal{N}(0, 1/s)$  random  
 116 variables. For a diagonal matrix  $R \in \mathbb{R}^{n \times n}$ , we call it a diagonal Rademacher matrix if its entries  
 117 are i.i.d. Rademacher random variables. See Appendix A for more detailed introductions.

### 118 4 Problem Setting

119 Throughout this paper, we consider a graph  $\mathcal{G}$  with  $n$  nodes and  $m$  edges. We denote its adjacency  
 120 matrix by  $A \in \mathbb{R}^{n \times n}$ , degree matrix by  $D = \text{diag}(A\mathbf{1}) \in \mathbb{R}^{n \times n}$ , where  $\mathbf{1} \in \mathbb{R}^n$  is a vector with  
 121 all entries equal to 1, and Laplace matrix (or Laplacian) by  $L = D - A \in \mathbb{R}^{n \times n}$ . We also use the  
 122 normalized adjacency  $\hat{A} = D^{-1/2}AD^{-1/2}$  and normalized Laplacian  $\hat{L} = I - \hat{A}$ . Notice that  
 123 we assume there’s a self-loop for each node to avoid 0-degree nodes, which is commonly ensured  
 124 in GNN implementations [22, 18], and therefore  $D^{-1/2}$  always exists. We also use  $B \in \mathbb{R}^{m \times n}$  to  
 125 represent the incidence matrix of graph  $\mathcal{G}$ , that is, if the  $i$ -th edge in  $\mathcal{G}$  is  $(u_i, v_i)$ , then the  $i$ -th row of  
 126  $B$  is  $\delta_{u_i} - \delta_{v_i}$  (the order of  $u_i$  and  $v_i$  is arbitrary). It is not hard to see that  $L = B^{\top}B$ .

127 We also assume that for each node  $u$  in the graph, there is a feature vector  $\mathbf{x}_u \in \mathbb{R}^d$ , and a label  
 128 vector<sup>1</sup>  $\mathbf{y}_u \in \mathbb{R}^c$  attached with it. Let  $X \in \mathbb{R}^{n \times d}$  and  $Y \in \mathbb{R}^{n \times c}$  be the stack of  $\mathbf{x}_i$ -s and  $\mathbf{y}_i$ -s. We  
 129 assume  $n \geq \max\{d, c\}$  and  $X$  and  $Y$  are both column full rank.

130 The Unfolded GNN we consider in this paper is based on TWIRLS [46], which is defined by  
 131 optimizing the following objective (called the “energy function”):

$$E(\mathbf{Z}; \mathbf{W}) := \frac{\lambda}{2} \text{Tr} \left( \mathbf{Z}^{\top} \hat{L} \mathbf{Z} \right) + \frac{1}{2} \|\mathbf{Z} - f(\mathbf{X}; \mathbf{W})\|_{\mathcal{F}}^2. \quad (4.1)$$

<sup>1</sup>We allow the label to be a vector to make the concept more general. When labels are one-hot vectors, the task is a classification task, and when labels are scalars, it is a regression task.

132 Here  $\mathbf{Z} \in \mathbb{R}^{n \times c}$  is the optimization variable that can be viewed as the hidden state of the model,  
 133  $\mathbf{W}$  is the learnable parameter of the model, and  $\lambda > 0$  is a hyper-parameter to balance the graph  
 134 structure information and node feature information. Note that  $f : \mathbb{R}^d \rightarrow \mathbb{R}^c$  is a function of  $\mathbf{X}$  and is  
 135 parameterized by matrix  $\mathbf{W}$ .

136 For the model training, given a loss function  $\ell : \mathbb{R}^{n \times c} \times \mathbb{R}^{n \times c} \rightarrow \mathbb{R}$ , the training of TWIRLS can be  
 137 represented by the following bi-level optimization problem

$$\mathbf{W}^* \in \arg \min_{\mathbf{W} \in \mathbb{R}^{d \times c}} \ell[\mathbf{Z}^*(\mathbf{W}); \mathbf{Y}] \quad (4.2)$$

$$\text{s.t. } \mathbf{Z}^*(\mathbf{W}) \in \arg \min_{\mathbf{Z} \in \mathbb{R}^{n \times c}} E(\mathbf{Z}; \mathbf{W}), \quad (4.3)$$

138 where the solution of inner problem eq.(4.3) (i.e.  $\mathbf{Z}^*(\mathbf{W})$ ) is viewed as the model output. The outer  
 139 problem eq.(4.2) is viewed as the problem of training.

140 In [46], the forward process of TWIRLS is obtained by unfolding the gradient descent algorithm for  
 141 minimizing  $E$  given a fixed  $\mathbf{W}$ :

$$\mathbf{Z}^{(t+1)} = \mathbf{Z}^{(t)} - \alpha \nabla E(\mathbf{Z}^{(t)}) \quad (4.4)$$

$$= (1 - \alpha\lambda - \alpha)\mathbf{Z}^{(t)} + \alpha \hat{\mathbf{A}} \mathbf{Z}^{(t)} + \alpha f(\mathbf{X}; \mathbf{W}), \quad (4.5)$$

142 where  $\alpha > 0$  is the step size of gradient descent. Notice that we fix  $\mathbf{W}$  and view  $E$  as a function of  
 143  $\mathbf{Z}$ . The model output is  $\mathbf{Z}^{(T)}$  for a fixed iteration number  $T$ . When  $\alpha$  is set to a suitable value and  $T$   
 144 is large, it is clear that  $\mathbf{Z}^{(T)}$  is an approximation of  $\mathbf{Z}^*$  defined in eq.(4.3).

145 When  $T \rightarrow \infty$ , the output of TWIRLS converges to the unique global optimum of eq.(4.1), which  
 146 has a closed form

$$\mathbf{Z}^*(\mathbf{W}) = \arg \min_{\mathbf{Z}} E(\mathbf{Z}) = \left( \mathbf{I} + \lambda \hat{\mathbf{L}} \right)^{-1} f(\mathbf{X}; \mathbf{W}). \quad (4.6)$$

147 In this paper, we always assume  $T \rightarrow \infty$  and thus we focus on the closed-form solution defined by  
 148 eq.(4.6).

149 The training problem is then defined as optimizing  $\mathbf{W}$  to minimize the loss function between the  
 150 model output  $\mathbf{Z}^*(\mathbf{W})$  and label  $\mathbf{Y}$ , i.e.  $\ell : \mathbb{R}^{n \times d} \times \mathbb{R}^{n \times c} \rightarrow \mathbb{R}$ .

## 151 5 Algorithm and Analysis

152 In this section we present our main algorithm HERTA, give the convergence result and analyze its  
 153 time complexity. We also introduce our key techniques and give a proof sketch for the main result.  
 154 To start we introduce the following notion.

155 **Definition 5.1** (Effective Laplacian dimension). *For normalized Laplacian  $\hat{\mathbf{L}}$  and regularization term*  
 156  $\lambda > 0$ , define  $n_\lambda := \text{Tr} \left[ \hat{\mathbf{L}} \left( \hat{\mathbf{L}} + \lambda^{-1} \mathbf{I} \right)^{-1} \right]$  to be the effective Laplacian dimension of the graph.

157 Denote  $\{\lambda_i\}_{i=1}^n$  as the eigenvalues of  $\hat{\mathbf{L}}$ . Then the effective Laplacian dimension can be written as  
 158  $n_\lambda = \sum_i \frac{\lambda_i}{\lambda_i + \lambda^{-1}}$ . We can see that  $n_\lambda$  is roughly the number of eigenvalues of  $\hat{\mathbf{L}}$  which are of order  
 159  $\Omega(\lambda^{-1})$ , i.e., the number of ‘‘large eigenvalues’’ mentioned in Theorem 1.1. It is not hard to see that  
 160  $n_\lambda \leq n$  and  $n_\lambda \rightarrow n$  as  $\lambda \rightarrow \infty$ .

161 Intuitively,  $n_\lambda$  represents the number of eigen-directions for which the Laplacian regularizer  
 162  $\frac{\lambda}{2} \text{Tr}(\mathbf{Z}^\top \hat{\mathbf{L}} \mathbf{Z})$  is large, thus having a significant impact on the energy function (4.1). For the  
 163 remaining eigen-directions which are less significant, the effect of the regularizer is minimal. This  
 164 can also be seen from the closed form solution of the model (4.6): when we apply matrix  $(\mathbf{I} + \lambda \hat{\mathbf{L}})^{-1}$   
 165 to a vector, the effect of the Laplacian is dominated by the  $\mathbf{I}$  term if this vector is aligned with a  
 166 ‘‘small’’ Laplacian eigen-direction (i.e., with an eigenvalue significantly smaller than  $\lambda^{-1}$ ), and thus  
 167 the objective is close to the usual unregularized least squares loss.

168 Based on the above discussion, it stands to reason that the larger the effective Laplacian dimension  
 169  $n_\lambda$  of the graph, the more work we must do during the training to preserve the graph structure of the  
 170 GNN.

171 **5.1 Main Result**

172 In this subsection, we present our main theoretical result. Before delving into the specific theorem,  
 173 we first outline the model implementation used in our theoretical analysis. Although we make  
 174 specific assumptions for this analysis, the experiments (see Section 6) demonstrate that our proposed  
 175 algorithm is effective in broader settings beyond those considered here.

176 Principally, the implementation of  $f(\mathbf{X}; \mathbf{W})$  can be arbitrary as long as it can be trained (i.e.  
 177 computable, smooth, etc.), however, we noticed that in [46], most of the SOTA results are achieved  
 178 by the implementation  $f(\mathbf{X}; \mathbf{W}) = \mathbf{X}\mathbf{W}$  where  $\mathbf{W} \in \mathbb{R}^{d \times c}$ . Therefore, hereinafter we focus on  
 179 this specific implementation.

180 For the simplicity of analysis, we present our theoretical result based on the mean squared error  
 181 (MSE) loss defined as  $\ell(\mathbf{Z}, \mathbf{Y}) := \frac{1}{2} \|\mathbf{Z} - \mathbf{Y}\|_{\mathcal{F}}^2$ . Despite this specific choice in the analysis, we  
 182 empirically show in Section 6 that our algorithm also works for cross entropy (CE) loss, which is the  
 183 most commonly used loss function for node classification tasks.

184 Under MSE loss, we can decompose the loss into sub-losses for each class, i.e.  $\ell[\mathbf{Z}^*(\mathbf{W}); \mathbf{Y}] =$   
 185  $\sum_{i=1}^c \ell_i(\mathbf{w}_i; \mathbf{y}_i)$ , where  $\mathbf{w}_i$  and  $\mathbf{y}_i$  are the  $i$ -th column of  $\mathbf{W}$  and  $\mathbf{Y}$  respectively. In the following  
 186 analysis, we will mostly focus on each  $\ell_i$ . Moreover, we fix  $\mathbf{y}_i$ , and view  $\ell_i$  as a function of  $\mathbf{w}_i$ .

187 **Theorem 5.1** (Main result). *For any  $\epsilon > 0$ , with a proper step size  $\eta$ , number of iterations  $T$  and*  
 188 *constant  $K > 0$ , Algorithm 1 finds a solution  $\hat{\mathbf{W}} \in \mathbb{R}^{d \times c}$  in time*

$$\tilde{O} \left( (m + nd) c \left( \log \frac{1}{\epsilon} \right)^2 + n_\lambda \lambda^2 d + d^3 \right) \quad (5.1)$$

189 *such that with probability  $1 - \frac{1}{n}$ ,  $\ell[\mathbf{Z}^*(\hat{\mathbf{W}}); \mathbf{Y}] \leq (1 + \epsilon) \cdot \ell^*$ , where  $\ell^* := \min_{\mathbf{W}} \ell[\mathbf{Z}^*(\mathbf{W}); \mathbf{Y}]$ .*

190 Notice that if we assume  $n_\lambda = O(n/\lambda^2)$ , then  $\tilde{O}(n_\lambda \lambda^2 d) = \tilde{O}(nd)$ . By further assuming the number  
 191 of classes  $c = O(1)$ , we recover the bound given by Theorem 1.1.

192 Before heading into the details of HERTA (Algorithm 1), in next subsection we first analyze the  
 193 original (standard) implementation of training TWIRLS and point out the efficiency bottlenecks. We  
 194 then introduce the key techniques and methods we need to prove our main result. Finally we present  
 195 our main algorithm HERTA, and give a proof sketch for Theorem 5.1.

196 **5.2 Analysis of TWIRLS Training**

197 In this subsection, we introduce the framework that we use to analyze the complexity of TWIRLS  
 198 training. We also provide a complexity analysis of the original implementation used in [46].

199 As noted in Section 4, the whole optimization problem can be viewed as a bi-level optimization prob-  
 200 lem, where the inner-problem eq.(4.3) approximates the linear system solver  $(\mathbf{I} + \lambda \hat{\mathbf{L}})^{-1} (\mathbf{X} \mathbf{w}_i -$   
 201  $\hat{\mathbf{y}}_i)$ , and the outer problem eq.(4.2) uses the linear system solution to approximate the optimal solution  
 202 of the training loss. When we discuss the solution of outer problems, we treat the solver of the inner  
 203 problem (linear solver) as a black box. We formally define a linear solver as follows.

204 **Definition 5.2** (Linear solver). *For a positive definite matrix  $\mathbf{H} \in \mathbb{R}^{n \times n}$  and a real number  $\epsilon > 0$ ,*  
 205 *we call  $f : \mathbb{R}^n \rightarrow \mathbb{R}^n$  a linear solver for  $\mathbf{H}$  with  $\epsilon$  error rate if it satisfies*

$$\forall \mathbf{u} \in \mathbb{R}^n, \|f(\mathbf{u}) - \mathbf{H}^{-1} \mathbf{u}\|_{\mathbf{H}} \leq \epsilon \|\mathbf{H}^{-1} \mathbf{u}\|_{\mathbf{H}}. \quad (5.2)$$

206 In [46], the outer problem is also solved by standard gradient descent. The gradient of  $\ell_i(\mathbf{w}_i)$  is

$$\nabla \ell_i(\mathbf{w}_i) = \mathbf{X}^\top \left( \mathbf{I} + \lambda \hat{\mathbf{L}} \right)^{-2} (\mathbf{X} \mathbf{w}_i - \hat{\mathbf{y}}_i), \quad (5.3)$$

207 where  $\hat{\mathbf{y}}_i = (\mathbf{I} + \lambda \hat{\mathbf{L}}) \mathbf{y}_i$ .

208 In actual implementation, the matrix inverse is calculated by linear solvers. Suppose that we have  
 209 a linear solver for  $(\mathbf{I} + \lambda \hat{\mathbf{L}})$  denoted by  $\mathcal{S}$ , then by embedding it into eq.(5.3) we can obtain the  
 210 gradient approximation as

$$\widetilde{\nabla \ell_i(\mathbf{w}_i)} := \mathbf{X}^\top \mathcal{S} [\mathcal{S} (\mathbf{X} \mathbf{w}_i - \hat{\mathbf{y}}_i)]. \quad (5.4)$$

211 The following convergence result can be derived for this approximate gradient descent. See Appendix  
 212 for the proof.

213 **Lemma 5.1** (Convergence). *If we minimize  $\ell_i$  using gradient descent with the gradient approx-*  
 214 *imation defined in eq.(5.4), where  $\mathcal{S}$  is a linear solver of  $\mathbf{I} + \lambda\hat{\mathbf{L}}$  with  $\mu$  error rate satisfying*  
 215  $\mu \leq \min \left\{ \frac{\epsilon^{1/2}}{50\kappa(\mathbf{X})\lambda^2}, 1 \right\}$ , *then, to obtain an  $\epsilon \in (0, 1)$  error rate the outer problem (i.e.  $\ell_i$ ), the*  
 216 *number of iterations needed is  $T = O\left(\kappa \log \frac{1}{\epsilon}\right)$ , where  $\kappa := \kappa(\mathbf{X}^\top (\mathbf{I} + \lambda\hat{\mathbf{L}})^{-2} \mathbf{X})$ .*

217 Based on Lemma 5.1, in Appendix B we provide a detailed analysis of the time complexity of  
 218 the implementation used by [46]. From this analysis, one can figure out two bottlenecks of the  
 219 computational complexity for solving TWIRLS: 1) The number of iterations needed for the inner  
 220 loop is dependent on  $\lambda$ ; 2) The number of iterations needed for the outer loop is dependent on the  
 221 condition number of the outer problem. Our acceleration algorithm is based on the following two  
 222 observations, which correspond to the aforementioned bottlenecks:

- 223 1. The data matrix in the inner loop is a graph Laplacian plus identity. It is possible to solve  
 224 the linear system defined by this matrix faster by exploiting its structure.
- 225 2. A good enough preconditioner can be constructed for the outer iteration using techniques  
 226 from RandNLA.

### 227 5.3 Key Techniques

228 In this subsection, we introduce the mathematical tools that we will use in the proof of our main  
 229 result.

230 **SDD Solvers.** If a symmetric matrix  $\mathbf{H} \in \mathbb{R}^{n \times n}$  satisfies

$$\forall 1 \leq k \leq n, \mathbf{H}_{k,k} \geq \sum_{i=1}^n |\mathbf{H}_{k,i}|, \quad (5.5)$$

231 it is called a symmetric diagonal dominated (SDD) matrix. It has been shown that the linear system  
 232 defined by sparse SDD matrices can be solved in an efficient way, via a SDD solver [21, 31].

233 Since  $\mathbf{I} + \lambda\hat{\mathbf{L}}$  is indeed a SDD matrix in our problem, we can apply off-the-shelf SDD solvers. In the  
 234 following, we use  $\text{SSolve}_\epsilon(\mathbf{H}, \mathbf{u})$  to denote the SDD solver that calculates  $\mathbf{H}^{-1}\mathbf{u}$  with error rate less  
 235 or equal to  $\epsilon^2$ . From Lemma A.1, the time complexity of calculating  $\text{SSolve}_\epsilon(\mathbf{H}; \mathbf{u})$  is  $\tilde{O}[\text{nnz}(\mathbf{H})]$ .  
 236 See Appendix A.2 for details.

237 **Fast Matrix Multiplication.** For a matrix  $\mathbf{Q} \in \mathbb{R}^{n \times d}$  where  $n > d$ , it is known that calculating  
 238  $\mathbf{Q}^\top \mathbf{Q}$  takes  $O(nd^2)$  time if implemented by brute force. In HERTA, we manage to achieve a faster  
 239 calculation of  $\mathbf{Q}^\top \mathbf{Q}$  through Subsampled Randomized Hadamard Transformation (SRHT) [37, 41].  
 240 The key idea of SRHT is to first apply a randomized Hadamard transformation to the matrix to make  
 241 the ‘‘information’’ evenly distributed in each row, then uniformly sample the rows of  $\mathbf{Q}$  to reduce  
 242 the dimension. In the following, we use Hadamard to represent the Hadamard transformation, see  
 243 Appendix A.3 for the definition of Hadamard as well as a detailed introduction.

### 244 5.4 Regularized Spectral Sparsifier

245 It has been shown that for any connected graph with  $n$  nodes, (no matter how dense it is) there is  
 246 always a sparsified graph with  $\tilde{O}\left(\frac{n}{\epsilon^2}\right)$  edges whose Laplacian is an  $\epsilon$ -spectral approximation of the  
 247 original graph Laplacian [36]. Although the existence is guaranteed, it is generally hard to construct  
 248 such a sparsified graph. An algorithm that finds such a sparsified graph is called a spectral sparsifier.

249 While there has been many existing explorations on constructing spectral sparsifier of a given graph  
 250 with high probability and tolerable time complexity [35, 36, 23], the setting considered in this paper  
 251 is different from the standard one: instead of graph Laplacian  $\mathbf{L}$ , the matrix we concern is the  
 252 normalized and regularized Laplacian  $\hat{\mathbf{L}} + \lambda^{-1}\mathbf{I}$ .

<sup>2</sup>Principally,  $\text{SSolve}(\mathbf{H}; \cdot)$  is a vector function, but for convenience we sometimes also apply it to matrices,  
 in which case we mean applying it column-wisely. We also use the same convention for other vector functions.

253 At first glance, a spectral sparsifier also works for our problem since  $\tilde{\mathbf{L}} \approx_\epsilon \hat{\mathbf{L}}$  implies  $\tilde{\mathbf{L}} + \lambda^{-1}\mathbf{I} \approx_\epsilon$   
 254  $\hat{\mathbf{L}} + \lambda^{-1}\mathbf{I}$ . However, it turns out that it is possible to make the bound even tighter if we take into  
 255 account the regularization term  $\lambda^{-1}\mathbf{I}$  here. A similar argument is also made in [4]. We note that the  
 256 setting explored in [4] is also different from ours: we consider the normalized Laplacian of a graph,  
 257 which is not a Laplacian of any graph (even allowing weighted edges).

258 To this end, we propose a new spectral sparsifier that works for  $\hat{\mathbf{L}} + \lambda^{-1}\mathbf{I}$  which is used in our  
 259 problem. This spectral sparsifier reduces the number of edges to  $O\left(\frac{n}{\epsilon^2} \log n\right)$ , which, as we discussed  
 260 before, is smaller than  $O\left(\frac{n}{\epsilon^2} \log n\right)$  obtained by standard spectral sparsifier.

261 **Lemma 5.2** (Regularized spectral sparsifier). *Let  $\hat{\mathbf{L}}$  be a (normalized) graph Laplacian with  $m$  edges*  
 262 *and  $n$  nodes. There exists a constant  $C$  such that for any  $\epsilon > 0$  and  $\lambda > 0$ , Algorithm 2 outputs*  
 263 *a (normalized) graph Laplacian  $\tilde{\mathbf{L}}$  with  $n$  nodes and  $O\left(\frac{n}{\epsilon^2} \log n\right)$  edges in time  $\tilde{O}(m)$ , such that*  
 264  *$\tilde{\mathbf{L}} + \lambda^{-1}\mathbf{I} \approx_\epsilon \hat{\mathbf{L}} + \lambda^{-1}\mathbf{I}$  with probability at least  $1 - \frac{1}{2n}$ .*

265 The basic idea behind Lemma 5.2 is to use ridge leverage score sampling for the normalized incidence  
 266 matrix  $\mathbf{B}\mathbf{D}^{-1/2}$  [10], which reduces the problem to estimating a regularized version of the effective  
 267 resistance. We modify the existing methods of estimating effective resistance [38] to make it work  
 268 for the normalized and regularized Laplacian.

---

**Algorithm 1** HERTA: A High-Efficiency and Rigorous Training Algorithm for TWIRLS

---

**Input:**  $\hat{\mathbf{L}}, \mathbf{X}, \mathbf{Y}, \lambda, c, K > 0$ , step size  $\eta$  and number of iteration  $T$ .  
 Set  $\beta = \frac{1}{64}$ ,  $s = \frac{Kd}{\beta^2} \log n$ ,  $\mu = \min\left\{\frac{\epsilon^{1/2}}{50\kappa(\mathbf{X})\lambda^2}, 1\right\}$   
 and  $\mathbf{R} \in \mathbb{R}^{n \times n}$  a diagonal Rademacher matrix;  
 $\tilde{\mathbf{L}} \leftarrow \text{Sparsify}_{\frac{\beta}{3\lambda}}(\hat{\mathbf{L}})$  by calling Algorithm 2;  
 $\mathbf{Q} \leftarrow \text{SSolve}_{\frac{\beta}{\sqrt{3\lambda}}}\left(\mathbf{I} + \lambda\tilde{\mathbf{L}}; \mathbf{X}\right)$ ;  
 $\mathbf{Q}' \leftarrow \text{Hadamard}(\mathbf{R}\mathbf{Q})$ ;  
 Subsample  $s$  rows of  $\mathbf{Q}'$  uniformly and obtain  $\tilde{\mathbf{Q}}$ ;  
 $\mathbf{P} \leftarrow \tilde{\mathbf{Q}}^\top \tilde{\mathbf{Q}}$ ;  
 $\mathbf{P}' \leftarrow \mathbf{P}^{-1/2}$ ;  
 269 **for**  $i = 1$  **to**  $c$  **do**  
    $\hat{\mathbf{y}}_i \leftarrow (\mathbf{I} + \lambda\tilde{\mathbf{L}})\mathbf{y}_i$ , where  $\mathbf{y}_i$  is the  $i$ -th column of  $\mathbf{Y}$ ;  
   Initialize  $w_i^{(0)}$  by all zeros;  
   **for**  $t = 1$  **to**  $T$  **do**  
      $\mathbf{u}_i^{(t)} \leftarrow \mathbf{X}\mathbf{P}'w_i^{(t-1)} - \hat{\mathbf{y}}_i$ ;  
      $\mathbf{u}_i^{(t)'} \leftarrow \text{SSolve}_\mu\left(\mathbf{I} + \lambda\tilde{\mathbf{L}}; \mathbf{u}_i^{(t)}\right)$ ;  
      $\mathbf{u}_i^{(t)''} \leftarrow \text{SSolve}_\mu\left(\mathbf{I} + \lambda\hat{\mathbf{L}}; \mathbf{u}_i^{(t)'}\right)$ ;  
      $\mathbf{g}_i^{(t)} \leftarrow \mathbf{P}'\mathbf{X}^\top \mathbf{u}_i^{(t)''}$ ;  
      $w_i^{(t)} \leftarrow w_i^{(t-1)} - \eta\mathbf{g}_i^{(t)}$ ;  
   **end for**  
**end for**  
**Output:**  $\left\{w_i^{(T)}\right\}_{i=1}^c$ .

---



---

**Algorithm 2** Sparsify $_\epsilon(\hat{\mathbf{L}})$ : Regularized Spectral Sparsifier

---

**Input:**  $\hat{\mathbf{L}}, \hat{\mathbf{B}}, \lambda, C > 0$  and expected error rate  $\epsilon$ .  
 Set  $k = C \log m$  and  $s = \frac{Cn_\lambda \log n}{\epsilon^2}$ , and construct Gaussian sketch  $\mathbf{\Pi}_1 \in \mathbb{R}^{k \times m}$ ,  $\mathbf{\Pi}_2 \in \mathbb{R}^{k \times n}$ ;  
 $\mathbf{B}_S \leftarrow \text{SSolve}_{2^{1/4}}(\hat{\mathbf{L}} + \lambda^{-1}\mathbf{I}; (\mathbf{\Pi}_1 \hat{\mathbf{B}})^\top)$ ;  
 $\mathbf{\Pi}_S \leftarrow \text{SSolve}_{2^{1/4}}(\hat{\mathbf{L}} + \lambda^{-1}\mathbf{I}; \mathbf{\Pi}_2^\top)$ ;  
**for**  $i = 1$  **to**  $m$  **do**  
    $\tilde{l}_i \leftarrow \|\mathbf{B}_S \hat{\mathbf{b}}_i\|^2 + \lambda^{-1}\|\mathbf{\Pi}_S \hat{\mathbf{b}}_i\|^2$ , where  $\hat{\mathbf{b}}_i^\top$  is the  $i$ -th row of  $\hat{\mathbf{B}}$ ;  
**end for**  
 $Z \leftarrow \sum_{i=1}^m \tilde{l}_i$ ;  
 Subsample  $s$  rows of  $\hat{\mathbf{B}}$  with probabilities  $\left\{\frac{\tilde{l}_i}{Z}\right\}_{i=1}^m$  and obtain  $\tilde{\mathbf{B}}$ ;  
 $\tilde{\mathbf{L}} \leftarrow \tilde{\mathbf{B}}^\top \tilde{\mathbf{B}}$ ;  
**Output:**  $\tilde{\mathbf{L}}$ .

---

## 270 5.5 Main Algorithm

271 Now we are ready to present our main algorithm HERTA, see Algorithm 1. As mentioned before,  
 272 HERTA is composed of two major components: constructing a preconditioner matrix  $\mathbf{P}$  and applying  
 273 it to the optimization problem. Below we introduce each part.

274 **Constructing the Preconditioner.** For the outer problem which can be ill-conditioned, we first  
 275 precondition it using a preconditioner  $\mathbf{P} \in \mathbb{R}^{d \times d}$  which is a constant level spectral approximation of

276 the Hessian, i.e.  $\mathbf{X}^\top (\mathbf{I} + \lambda \hat{\mathbf{L}})^{-2} \mathbf{X}$ . We claim that the matrix  $\mathbf{P}$  constructed inside Algorithm 1  
 277 indeed satisfies this property.

278 **Lemma 5.3** (Preconditioner). *Let  $\mathbf{P}$  be the matrix constructed in Algorithm 1. We have  $\mathbf{P} \approx_{\frac{1}{2}}$   
 279  $\mathbf{X}^\top (\mathbf{I} + \lambda \hat{\mathbf{L}})^{-2} \mathbf{X}$  with probability at least  $1 - \frac{1}{n}$ .*

280 **Solving the Outer Problem.** After obtaining the preconditioner  $\mathbf{P}$  which approximates the Hessian  
 281 by a constant level, we use it to precondition the outer problem and provably reduce the condition  
 282 number to a constant. With the new problem being well-conditioned, iterative methods (like gradient  
 283 descent) take much less steps to converge.

284 **Lemma 5.4** (Well-conditioned Hessian). *Let  $\ell'$  be such that  $\ell'(\mathbf{w}_i) = \ell(\mathbf{P}^{-1/2} \mathbf{w}_i)$ . Suppose for  
 285 some constant  $c_0 \in (0, 1)$  we have  $\mathbf{P} \approx_{c_0} \mathbf{X}^\top (\mathbf{I} + \lambda \hat{\mathbf{L}})^{-2} \mathbf{X}$ , then the condition number of the  
 286 Hessian of  $\ell'$  is bounded by*

$$\kappa(\nabla^2 \ell') \leq (1 + c_0)^2. \quad (5.6)$$

287 *Moreover, if  $\mathbf{w}'$  is a solution of  $\ell'$  with  $\epsilon$  error rate, then  $\mathbf{P}^{-1/2} \mathbf{w}'$  is a solution of  $\ell$  with  $\epsilon$  error rate.*

288 Notice that, the problem  $\ell'$  defined in Lemma 5.4 can be viewed as the original problem  $\ell$  with  $\mathbf{X}$   
 289 being replaced by  $\mathbf{X} \mathbf{P}^{-1/2}$ , therefore we can use Lemma 5.1 and obtain the convergence rate of  
 290 HERTA. With the convergence result and analysis of the running time for each step, we are able to  
 291 prove Theorem 5.1. See Appendix C.7 for the full proof.

292 **Remark on the Unavoidable  $\lambda^2$  in Runtime.** Notice that  $\lambda^2$  appears in the time bound given in  
 293 Theorem 5.1. From the proof in Appendix C.5 we can see that this term originates from the squared  
 294 inverse of  $\mathbf{I} + \lambda \hat{\mathbf{L}}$ . In Appendix D.2, we show that in the worst case even if we approximate a matrix  
 295 to a very high precision, the approximation rate of the corresponding squared version can still be  
 296 worse by a factor of the condition number of the matrix, which in our case leads to the unavoidable  
 297  $\lambda^2$  in the runtime.

298 **Remark on Sparsifying the Graph in Each Iteration.** Note that in Algorithm 1, we use the  
 299 complete Laplacian  $\hat{\mathbf{L}}$  for gradient iterations (i.e., the for-loop), and the graph sampling only occurs  
 300 when constructing the preconditioner  $\mathbf{P}$ . We note that sampling  $\hat{\mathbf{L}}$  in the for-loop can lead to extra  
 301 loss in the gradient estimation, which forces us to sparsify the graph to a very high precision to  
 302 ensure an accurate enough gradient for convergence. This could result in a suboptimal running  
 303 time. Moreover, as discussed in Section 1, the current running time bound for HERTA is optimal  
 304 up to logarithmic factors, which means that there is little to gain from performing extra sampling in  
 305 for-loops. See Appendix D.3 for a more detailed and quantitative discussion.

## 306 6 Experiments

307 In this section, we verify our theoretical results through experiments on real world datasets. Since for  
 308 the inner problem we use off-the-shelf SDD solvers, in this section we focus on the outer problem,  
 309 i.e. the training loss convergence rate. In each setting we compare the training loss convergence rate  
 310 of TWIRLS trained by our method against that trained by standard gradient descent using exactly the  
 311 same training hyper-parameters.

312 **Datasets.** We conduct experiments on all of the datasets used in the Section 5.1 of [46]: Cora,  
 313 Citeseer and Pubmed collected by [34], as well as the ogbn-arxiv dataset from the Open Large Graph  
 314 Benchmark (OGB, 19).

315 Notice that for Cora, Citeseer and Pubmed, common practice uses the semi-supervised setting  
 316 where there is a small training set and relatively large validation set. As we are comparing training  
 317 convergence rate, we find it more comparative to use a larger training set (which makes solving the  
 318 optimization problem of training more difficult). Therefore, we randomly select 80% of nodes as  
 319 training set for Cora, Citeseer and Pubmed. For OGB, we use the standard training split.

320 Due to limited space, we only include the results with  $\lambda = 1$  and with datasets Citeseer, Pubmed and  
 321 ogbn-arxiv in this section, and delay additional results to Appendix F.

322 **6.1 Convergence Rate Comparison Under MSE Loss**

323 In this section, we compare the convergence rate for models trained with MSE loss, which is well  
 324 aligned with the setting used to derive our theory. We also adopt a variation of HERTA that allows  
 325 using it on other optimizers and apply it on Adam optimizer. See Figure 1 for the results. Notice  
 326 that for each figure, we shift the curve by the minimum value of the loss (which is approximated by  
 327 running HERTA for more epochs) and take logarithm scale to make the comparison clearer.

328 From the results, it is clear that on all datasets and all optimizers we consider, HERTA converges  
 329 much faster than standard methods. It generally requires  $\leq 10$  iterations to get very close to the  
 330 smallest training loss. This shows that the guarantee obtained in our theoretical results (Theorem 5.1)  
 331 not only matches the experiments, but also holds when the setting is slightly changed (using other  
 332 optimizers). Thus, these experimental results verify the universality of HERTA as an optimization  
 333 framework.

334 We also conduct extensive experiments with larger  $\lambda$ . See Appendix F for the results. These results  
 335 verify that even when  $\lambda$  is relatively large, HERTA converges very fast, which suggests that the  
 336 dominant term in Theorem 5.1 is the  $\tilde{O}(m)$  term instead of the  $\tilde{O}(n_\lambda \lambda^2 d)$  term.

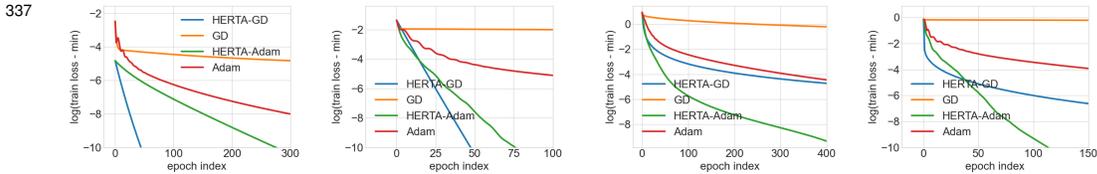


Figure 1: The training loss comparison between HERTA and standard optimizers on MSE loss with  $\lambda = 1$ . Left: ogbn-arxiv; Right: pubmed.

Figure 2: The training loss comparison between HERTA and standard optimizers on CE loss with  $\lambda = 1$ . Left: ogbn-arxiv; Right: pubmed.

338 **6.2 Convergence Rate Comparison under Cross Entropy Loss**

339 Note that in the theoretical analysis of Section 5, we focus on the MSE loss. However, for graph  
 340 node classification tasks, MSE is not the most commonly used loss. Instead, most node classification  
 341 models use the cross entropy loss. For example, the original TWIRLS paper [46] uses CE loss.

342 To this end, we also adopt HERTA on training problems with CE loss. See Appendix E for the details  
 343 of the implementation. The results are displayed in Figure 2. From these results, it is clear that  
 344 HERTA also significantly speeds up training with CE loss. The results demonstrate that although  
 345 originally developed on MSE loss, HERTA is not limited to it and have the flexibility to be applied to  
 346 other type of loss functions.

347 **Remark on the Surprising Effectiveness of HERTA on Cross Entropy Loss.** In Section 6.2, we  
 348 observe that despite not being guaranteed by the theoretical result, HERTA shows a certain degree of  
 349 universality in that it also works on CE loss. We claim that this phenomenon might originate from the  
 350 fact the Hessians of TWIRLS under MSE loss and CE loss are very similar. We provide an analysis  
 351 of the gradient and Hessian of TWIRLS under MSE and CE losses in Appendix D.1. The results  
 352 show that the gradient under CE loss can be viewed as the gradient under MSE loss with one term  
 353 being normalized by softmax, and the Hessian under CE loss can be viewed as a rescaled version of  
 354 the Hessian under MSE loss. These comparisons serve as an explanation of why HERTA works so  
 355 well with CE loss.

356 **7 Conclusions**

357 In this paper we present HERTA: a High-Efficiency and Rigorous Training Algorithm that solves the  
 358 problem of training Unfolded GNNs on a graph within a time essentially of the same magnitude as  
 359 the time it takes to load the input data. As a component of HERTA, we also propose a new spectral  
 360 sparsifier that works for normalized and regularized graph Laplacian matrices.

361 Experimental results on real world datasets show the effectiveness of HERTA. Moreover, it is shown  
 362 that HERTA works for various loss functions and optimizers, despite being derived from a specific  
 363 loss function and optimizer. This shows the universality of HERTA and verifies that it is ready to use  
 364 in practice. See Appendix D for further discussions and possible extensions of HERTA.

## References

- [1] Hongjoon Ahn, Yongyi Yang, Quan Gan, Taesup Moon, and David P Wipf. Descent steps of a relation-aware energy produce heterogeneous graph neural networks. *Advances in Neural Information Processing Systems*, 35:38436–38448, 2022.
- [2] Peter W Battaglia, Jessica B Hamrick, Victor Bapst, Alvaro Sanchez-Gonzalez, Vinicius Zambaldi, Mateusz Malinowski, Andrea Tacchetti, David Raposo, Adam Santoro, Ryan Faulkner, et al. Relational inductive biases, deep learning, and graph networks. *arXiv preprint arXiv:1806.01261*, 2018.
- [3] Léon Bottou, Frank E. Curtis, and Jorge Nocedal. Optimization methods for large-scale machine learning. *SIAM Rev.*, 60(2):223–311, 2018.
- [4] Daniele Calandriello, Alessandro Lazaric, Ioannis Koutis, and Michal Valko. Improved large-scale graph learning through ridge spectral sparsification. In *International Conference on Machine Learning*, pages 688–697. PMLR, 2018.
- [5] Jianfei Chen, Jun Zhu, and Le Song. Stochastic training of graph convolutional networks with variance reduction. In Jennifer G. Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholmsmässan, Stockholm, Sweden, July 10-15, 2018*, volume 80 of *Proceedings of Machine Learning Research*, pages 941–949. PMLR, 2018.
- [6] Jie Chen, Tengfei Ma, and Cao Xiao. Fastgcn: fast learning with graph convolutional networks via importance sampling. *arXiv preprint arXiv:1801.10247*, 2018.
- [7] Wei-Lin Chiang, Xuanqing Liu, Si Si, Yang Li, Samy Bengio, and Cho-Jui Hsieh. Cluster-gcn: An efficient algorithm for training deep and large graph convolutional networks. In *Proceedings of the 25th ACM SIGKDD international conference on knowledge discovery & data mining*, pages 257–266, 2019.
- [8] Kenneth L Clarkson and David P Woodruff. Low-rank approximation and regression in input sparsity time. *Journal of the ACM (JACM)*, 63(6):1–45, 2017.
- [9] Michael B Cohen, Yin Tat Lee, and Zhao Song. Solving linear programs in the current matrix multiplication time. *Journal of the ACM (JACM)*, 68(1):1–39, 2021.
- [10] Michael B Cohen, Cameron Musco, and Christopher Musco. Input sparsity time low-rank approximation via ridge leverage score sampling. In *Proceedings of the Twenty-Eighth Annual ACM-SIAM Symposium on Discrete Algorithms*, pages 1758–1777. SIAM, 2017.
- [11] Michael B Cohen, Cameron Musco, and Christopher Musco. Input sparsity time low-rank approximation via ridge leverage score sampling. In *Proceedings of the Twenty-Eighth Annual ACM-SIAM Symposium on Discrete Algorithms*, pages 1758–1777. SIAM, 2017.
- [12] Sanjoy Dasgupta and Anupam Gupta. An elementary proof of a theorem of johnson and lindenstrauss. *Random Structures & Algorithms*, 22(1):60–65, 2003.
- [13] Michał Dereziński and Jiaming Yang. Solving dense linear systems faster than via preconditioning. *arXiv preprint arXiv:2312.08893*, 2023.
- [14] Mucong Ding, Kezhi Kong, Jingling Li, Chen Zhu, John Dickerson, Furong Huang, and Tom Goldstein. VQ-GNN: A universal framework to scale up graph neural networks using vector quantization. In Marc Aurelio Ranzato, Alina Beygelzimer, Yann N. Dauphin, Percy Liang, and Jennifer Wortman Vaughan, editors, *Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*, pages 6733–6746, 2021.
- [15] Petros Drineas and Michael W Mahoney. Randnla: randomized numerical linear algebra. *Communications of the ACM*, 59(6):80–90, 2016.

- 411 [16] Matthias Fey, Jan Eric Lenssen, Frank Weichert, and Jure Leskovec. Gnnautoscale: Scalable and  
412 expressive graph neural networks via historical embeddings. In Marina Meila and Tong Zhang,  
413 editors, *Proceedings of the 38th International Conference on Machine Learning, ICML 2021,*  
414 *18-24 July 2021, Virtual Event*, volume 139 of *Proceedings of Machine Learning Research*,  
415 pages 3294–3304. PMLR, 2021.
- 416 [17] Nathan Halko, Per-Gunnar Martinsson, and Joel A Tropp. Finding structure with randomness:  
417 Probabilistic algorithms for constructing approximate matrix decompositions. *SIAM review*,  
418 53(2):217–288, 2011.
- 419 [18] William L. Hamilton, Zhitao Ying, and Jure Leskovec. Inductive representation learning on  
420 large graphs. In Isabelle Guyon, Ulrike von Luxburg, Samy Bengio, Hanna M. Wallach, Rob  
421 Fergus, S. V. N. Vishwanathan, and Roman Garnett, editors, *Advances in Neural Information*  
422 *Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017,*  
423 *December 4-9, 2017, Long Beach, CA, USA*, pages 1024–1034, 2017.
- 424 [19] Weihua Hu, Matthias Fey, Marinka Zitnik, Yuxiao Dong, Hongyu Ren, Bowen Liu, Michele  
425 Catasta, and Jure Leskovec. Open graph benchmark: Datasets for machine learning on graphs.  
426 *Advances in neural information processing systems*, 33:22118–22133, 2020.
- 427 [20] Haitian Jiang, Renjie Liu, Xiao Yan, Zhenkun Cai, Minjie Wang, and David Wipf.  
428 Musegnn: Interpretable and convergent graph neural network layers at scale. *arXiv preprint*  
429 *arXiv:2310.12457*, 2023.
- 430 [21] Jonathan A Kelner, Lorenzo Orecchia, Aaron Sidford, and Zeyuan Allen Zhu. A simple,  
431 combinatorial algorithm for solving sdd systems in nearly-linear time. In *Proceedings of the*  
432 *forty-fifth annual ACM symposium on Theory of computing*, pages 911–920, 2013.
- 433 [22] Thomas N. Kipf and Max Welling. Semi-supervised classification with graph convolutional  
434 networks. *CoRR*, abs/1609.02907, 2016.
- 435 [23] Ioannis Koutis, Alex Levin, and Richard Peng. Faster spectral sparsification and numerical  
436 algorithms for SDD matrices. *ACM Trans. Algorithms*, 12(2):17:1–17:16, 2016.
- 437 [24] Meng Liu, Hongyang Gao, and Shuiwang Ji. Towards deeper graph neural networks. In  
438 *Proceedings of the 26th ACM SIGKDD international conference on knowledge discovery &*  
439 *data mining*, pages 338–348, 2020.
- 440 [25] Xiaorui Liu, Wei Jin, Yao Ma, Yaxin Li, Hua Liu, Yiqi Wang, Ming Yan, and Jiliang Tang.  
441 Elastic graph neural networks. In *International Conference on Machine Learning*, pages  
442 6837–6849. PMLR, 2021.
- 443 [26] Yao Ma, Xiaorui Liu, Tong Zhao, Yozen Liu, Jiliang Tang, and Neil Shah. A unified view on  
444 graph neural networks as graph signal denoising. In *Proceedings of the 30th ACM International*  
445 *Conference on Information & Knowledge Management*, pages 1202–1211, 2021.
- 446 [27] Per-Gunnar Martinsson and Joel A Tropp. Randomized numerical linear algebra: Foundations  
447 and algorithms. *Acta Numerica*, 29:403–572, 2020.
- 448 [28] Xiangrui Meng, Michael A Saunders, and Michael W Mahoney. Lsrn: A parallel iterative  
449 solver for strongly over-or underdetermined systems. *SIAM Journal on Scientific Computing*,  
450 36(2):C95–C118, 2014.
- 451 [29] Riley Murray, James Demmel, Michael W Mahoney, N Benjamin Erichson, Maksim Mel-  
452 nichenko, Osman Asif Malik, Laura Grigori, Piotr Luszczek, Michał Dereziński, Miles E Lopes,  
453 et al. Randomized numerical linear algebra: A perspective on the field with an eye to software.  
454 *arXiv preprint arXiv:2302.11474*, 2023.
- 455 [30] Kenta Oono and Taiji Suzuki. Graph neural networks exponentially lose expressive power for  
456 node classification. *arXiv preprint arXiv:1905.10947*, 2019.
- 457 [31] Richard Peng and Daniel A. Spielman. An efficient parallel solver for SDD linear systems. In  
458 David B. Shmoys, editor, *Symposium on Theory of Computing, STOC 2014, New York, NY, USA,*  
459 *May 31 - June 03, 2014*, pages 333–342. ACM, 2014.

- 460 [32] Richard Peng and Santosh Vempala. Solving sparse linear systems faster than matrix multiplication. In *Proceedings of the 2021 ACM-SIAM symposium on discrete algorithms (SODA)*, pages 461 504–521. SIAM, 2021.
- 462
- 463 [33] Vladimir Rokhlin and Mark Tygert. A fast randomized algorithm for overdetermined linear least- 464 squares regression. *Proceedings of the National Academy of Sciences*, 105(36):13212–13217, 465 2008.
- 466 [34] Prithviraj Sen, Galileo Namata, Mustafa Bilgic, Lise Getoor, Brian Galligher, and Tina Eliassi- 467 Rad. Collective classification in network data. *AI magazine*, 29(3):93–93, 2008.
- 468 [35] Daniel A Spielman and Nikhil Srivastava. Graph sparsification by effective resistances. In 469 *Proceedings of the fortieth annual ACM symposium on Theory of computing*, pages 563–568, 470 2008.
- 471 [36] Daniel A Spielman and Shang-Hua Teng. Spectral sparsification of graphs. *SIAM Journal on 472 Computing*, 40(4):981–1025, 2011.
- 473 [37] Joel A Tropp. Improved analysis of the subsampled randomized hadamard transform. *Advances 474 in Adaptive Data Analysis*, 3(01n02):115–126, 2011.
- 475 [38] Nisheeth K Vishnoi et al.  $Lx = b$  laplacian solvers and their algorithmic applications. *Founda- 476 tions and Trends® in Theoretical Computer Science*, 8(1–2):1–141, 2013.
- 477 [39] Nisheeth K Vishnoi et al.  $Lx = b$ . *Foundations and Trends® in Theoretical Computer Science*, 478 8(1–2):1–141, 2013.
- 479 [40] Minjie Wang, Da Zheng, Zihao Ye, Quan Gan, Mufei Li, Xiang Song, Jinjing Zhou, Chao Ma, 480 Lingfan Yu, Yu Gai, Tianjun Xiao, Tong He, George Karypis, Jinyang Li, and Zheng Zhang. 481 Deep graph library: A graph-centric, highly-performant package for graph neural networks. 482 *arXiv preprint arXiv:1909.01315*, 2019.
- 483 [41] David P Woodruff et al. Sketching as a tool for numerical linear algebra. *Foundations and 484 Trends® in Theoretical Computer Science*, 10(1–2):1–157, 2014.
- 485 [42] Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How powerful are graph neural 486 networks? *arXiv preprint arXiv:1810.00826*, 2018.
- 487 [43] Rui Xue, Haoyu Han, MohamadAli Torkamani, Jian Pei, and Xiaorui Liu. Lazygcn: Large-scale 488 graph neural networks via lazy propagation. *arXiv preprint arXiv:2302.01503*, 2023.
- 489 [44] Chaoqi Yang, Ruijie Wang, Shuochao Yao, Shengzhong Liu, and Tarek Abdelzaher. Revisiting 490 over-smoothing in deep gcns. *arXiv preprint arXiv:2003.13663*, 2020.
- 491 [45] Yongyi Yang, Tang Liu, Yangkun Wang, Zengfeng Huang, and David Wipf. Implicit vs unfolded 492 graph neural networks. *arXiv preprint arXiv:2111.06592*, 2021.
- 493 [46] Yongyi Yang, Tang Liu, Yangkun Wang, Jinjing Zhou, Quan Gan, Zhewei Wei, Zheng Zhang, 494 Zengfeng Huang, and David Wipf. Graph neural networks inspired by classical iterative 495 algorithms. In Marina Meila and Tong Zhang, editors, *Proceedings of the 38th International 496 Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event*, volume 139 of 497 *Proceedings of Machine Learning Research*, pages 11773–11783. PMLR, 2021.
- 498 [47] Hanqing Zeng, Muhan Zhang, Yinglong Xia, Ajitesh Srivastava, Andrey Malevich, Rajgopal 499 Kannan, Viktor Prasanna, Long Jin, and Ren Chen. Decoupling the depth and scope of graph 500 neural networks. *Advances in Neural Information Processing Systems*, 34:19665–19679, 2021.
- 501 [48] Hongwei Zhang, Tijin Yan, Zenjun Xie, Yuanqing Xia, and Yuan Zhang. Revisiting graph 502 convolutional network on semi-supervised node classification from an optimization perspective. 503 *arXiv preprint arXiv:2009.11469*, 2020.
- 504 [49] Jiong Zhu, Yujun Yan, Lingxiao Zhao, Mark Heimann, Leman Akoglu, and Danai Koutra. Be- 505 yond homophily in graph neural networks: Current limitations and effective designs. *Advances 506 in neural information processing systems*, 33:7793–7804, 2020.

- 507 [50] Meiqi Zhu, Xiao Wang, Chuan Shi, Houye Ji, and Peng Cui. Interpreting and unifying graph  
508 neural networks with an optimization framework. In *Proceedings of the Web Conference 2021*,  
509 pages 1215–1226, 2021.
- 510 [51] Daniel Zügner, Amir Akbarnejad, and Stephan Günnemann. Adversarial attacks on neural  
511 networks for graph data. In *Proceedings of the 24th ACM SIGKDD international conference on*  
512 *knowledge discovery & data mining*, pages 2847–2856, 2018.

513 **A Introduction to the Mathematical Tools**

514 In the main paper, due to space limitations, our introduction of certain mathematical tools and concepts  
 515 remains brief. In this section, we provide a more comprehensive introduction to the mathematical  
 516 tools utilized in the main paper.

517 **A.1 Subsampling**

518 Subsampling is used in multiple places in our Algorithm. When we say subsampling  $s$  rows of a  
 519 matrix  $M \in \mathbb{R}^{n \times d}$  with probability  $\{p_k\}_{k=1}^n$  and obtain a new matrix  $\tilde{M} \in \mathbb{R}^{s \times d}$ , it means each  
 520 row of  $\tilde{M}$  is an i.i.d. random vector which follows the same distribution of a random vector  $\xi$  that  
 521 satisfies the following property:

$$\mathbb{P} \left\{ \xi = \sqrt{\frac{1}{sp_k}} \mathbf{m}_k \right\} = p_k, \quad (\text{A.1})$$

522 where  $\mathbf{m}_k$  is the  $k$ -th row of matrix  $M$ . It is clear that constructing such a subsampled matrix  $\tilde{M}$   
 523 takes time  $O(sd)$ . Moreover, if  $p_k = \frac{1}{n}$ , we say the subsampling is with uniform probability.

524 The operator of subsampling  $s$  rows of a  $n$ -row matrix is clearly a linear transformation. Therefore  
 525 it can be represented by a random matrix  $S \in \mathbb{R}^{s \times n}$ . Notice that when we apply  $S$  to a matrix (i.e.  
 526  $S \cdot M$ ) we don't actually need to construct such a matrix  $S$  and calculate matrix product, since we  
 527 only carry out sampling.

528 **A.2 SDD Solvers**

529 As mentioned in the main paper, for a sparse SDD matrix  $M$ , it is possible to fast approximate its  
 530 linear solver. The following Lemma A.1 rigorously states the bound we can obtain for SDD solvers.

531 **Lemma A.1** ([31]). *For any SDD matrix  $M \in \mathbb{R}^{n \times n}$  and any real number  $\epsilon > 0$ , there exists an*  
 532 *algorithm to construct an operator  $\text{SSolve}_\epsilon(M; \cdot)$ , such that  $\text{SSolve}_\epsilon(M; \cdot)$  is a linear solver for*  
 533  *$M$  with  $\epsilon$  error rate and for any  $\mathbf{x} \in \mathbb{R}^n$ , calculating  $\text{SSolve}_\epsilon(M; \mathbf{x})$  takes time*

$$O \left( m \left( \log \frac{1}{\epsilon} \right) \text{poly log}(m) \text{poly log}(\kappa(M)) \right), \quad (\text{A.2})$$

534 where  $m = \text{nnz}(M)$ .

535 **A.3 Fast Matrix Multiplication**

536 As mentioned in the main paper, we use SRHT to achieve fast matrix multiplication. Here we further  
 537 explain this process. We refer interested readers to [37] for more details. For a positive integer  
 538 number  $n$  that is a power of 2, the Hadamard transformation of size  $n$  is a linear transformation  
 539 recursively defined as follows:

$$\mathbf{H}_n := \begin{bmatrix} \mathbf{H}_{n/2} & \mathbf{H}_{n/2} \\ \mathbf{H}_{n/2} & -\mathbf{H}_{n/2} \end{bmatrix} \quad (\text{A.3})$$

540 and  $\mathbf{H}_1 = 1$ . For a vector  $\mathbf{x} \in \mathbb{R}^n$ , we define

$$\text{Hadamard}(\mathbf{x}) := \frac{1}{\sqrt{n}} \mathbf{H}_n \mathbf{x}. \quad (\text{A.4})$$

541 Notice that when  $n$  is not a power of 2, we pad the vector  $\mathbf{x}$  with 0-s beforehand, therefore the  
 542 requirement of  $n$  being a power of 2 is ignored in the actual usage. From its recursive nature, it is not  
 543 hard to see that applying a Hadamard transformation to each column of an  $n \times d$  matrix only takes  
 544  $O(d \log n)$  time.

545 Let  $\mathbf{R} \in \mathbb{R}^{n \times n}$  be a diagonal matrix whose diagonal entries are i.i.d. Rademacher variables, and  
 546  $\mathbf{S} \in \mathbb{R}^{s \times n}$  be a subsampling matrix with uniform probability and  $s = O\left(\frac{d}{\beta^2} \log n\right)$ . For a matrix

547  $\mathbf{Q} \in \mathbb{R}^{n \times d}$ , its SRHT with  $\beta$  error rate is defined as

$$\text{SRHT}_\beta(\mathbf{Q}) := \mathbf{S} \text{Hadamard}(\mathbf{R}\mathbf{Q}), \quad (\text{A.5})$$

548 which is exactly the matrix  $\tilde{\mathbf{Q}}$  we use in Algorithm 1. A fast matrix multiplication result can be  
 549 achieved using SRHT.

550 **Lemma A.2** ([37]). For matrix  $\mathbf{Q} \in \mathbb{R}^{n \times d}$  where  $n \geq d$  and  $\text{rank}(\mathbf{Q}) = d$ , let  $\tilde{\mathbf{Q}} = \text{SRHT}_\beta(\mathbf{Q})$   
 551 where  $\beta \in (0, 1/4)$ , then we have

$$\tilde{\mathbf{Q}}^\top \tilde{\mathbf{Q}} \approx_\beta \mathbf{Q}^\top \mathbf{Q} \quad (\text{A.6})$$

552 with probability at least  $1 - \frac{1}{2^n}$ .

## 553 B Running Time of the Original Implementation Used in [46]

554 In this section we analyze the time complexity of the implementation used in [46], which uses  
 555 gradient descent to solve both inner and outer problem.

556 **Inner Problem Analysis.** We first analyze the time complexity of the inner problem solver used in  
 557 [46]. Recall that for the inner problem eq.(4.3), we need to find an approximation of the linear solver  
 558 of  $\mathbf{I} + \lambda \hat{\mathbf{L}}$ . In [46], it is implemented by a standard gradient descent. Here we consider approximately  
 559 solving the following least square problem by gradient descent:

$$\mathbf{v} = \arg \min_{\mathbf{v} \in \mathbb{R}^n} \frac{1}{2} \left\| (\mathbf{I} + \lambda \hat{\mathbf{L}}) \mathbf{v} - \mathbf{u} \right\|_2^2. \quad (\text{B.1})$$

560 The gradient step with step size  $\mu$  is:

$$\mathbf{v}^{(t+1)} = (1 - \mu) \left( \mathbf{I} + \lambda \hat{\mathbf{L}} \right)^2 \mathbf{v}^{(t)} + \eta \left( \mathbf{I} + \lambda \hat{\mathbf{L}} \right) \mathbf{u}. \quad (\text{B.2})$$

561 **Theorem B.1** (Inner Analysis). If we update  $\mathbf{v}^{(t)}$  through eq.(B.2) with initialization  $\mathbf{v}^{(0)} = \mathbf{0}$  and a  
 562 proper  $\eta$ , then after  $T = O\left(\lambda^2 \log \frac{\lambda}{\epsilon^2}\right)$  iterations we can get

$$\left\| \mathbf{v}^{(T)} - \left( \mathbf{I} + \lambda \hat{\mathbf{L}} \right)^{-1} \mathbf{u} \right\|_{(\mathbf{I} + \lambda \hat{\mathbf{L}})} \leq \epsilon \left\| \left( \mathbf{I} + \lambda \hat{\mathbf{L}} \right)^{-1} \mathbf{u} \right\|_{(\mathbf{I} + \lambda \hat{\mathbf{L}})} \quad (\text{B.3})$$

563 for any  $\epsilon \in (0, 1)$ .

564 *Proof.* Let  $\mathbf{H} = \mathbf{I} + \lambda \hat{\mathbf{L}}$ . Notice that the strongly convexity and Lipschitz constant of Problem  
 565 eq.(B.1) is  $\mathcal{C} = 1$  and  $L \leq 9\lambda^2$  respectively. From Lemma C.1, take  $\eta = \frac{1}{L}$  and  $\mathbf{v}^{(0)} = \mathbf{0}$ , we have

$$\left\| \mathbf{H} \mathbf{v}^{(T)} - \mathbf{u} \right\|^2 \leq (1 - 9\lambda^2)^T \left\| \mathbf{u} \right\|^2 \quad (\text{B.4})$$

566 for any  $T \in \mathbb{N}$ . Therefore,

$$\frac{\left\| \mathbf{v}^{(T)} - \mathbf{H}^{-1} \mathbf{u} \right\|_{\mathbf{H}}^2}{\left\| \mathbf{H}^{-1} \mathbf{u} \right\|_{\mathbf{H}}^2} = \frac{\left\| \mathbf{H} \mathbf{v}^{(T)} - \mathbf{u} \right\|_{\mathbf{H}^{-1}}^2}{\left\| \mathbf{u} \right\|_{\mathbf{H}^{-1}}^2} \quad (\text{B.5})$$

$$\leq 3\lambda \frac{\left\| \mathbf{H} \mathbf{v}^{(T)} - \mathbf{u} \right\|_2^2}{\left\| \mathbf{u} \right\|_2^2} \quad (\text{B.6})$$

$$\leq 3\lambda (1 - 9\lambda^2)^T. \quad (\text{B.7})$$

567 To obtain an error rate  $\epsilon$ , the smallest number of iterations  $T$  needed is

$$T = \left( \log \frac{1}{1 - (9\lambda^2)^{-1}} \right)^{-1} \log \frac{3\lambda}{\epsilon^2} + 1 = O \left[ \left( \log \frac{1}{1 - \lambda^{-2}} \right)^{-1} \log \frac{\kappa}{\epsilon^2} \right] = O \left( \lambda^2 \log \frac{\lambda}{\epsilon^2} \right). \quad (\text{B.8})$$

568  $\square$

569 **Overall Running Time.** As we proved above in Theorem B.1, the number of iterations needed for  
 570 solving the inner problem to error rate  $\beta$  with the implementation of [46] is  $\tilde{O}\left(\lambda^2 \log \frac{1}{\beta}\right)$ , which  
 571 is related to the hyper-parameter  $\lambda$ . For each inner iteration (i.e. eq.(B.2)), we need to compute  
 572  $(\mathbf{I} + \lambda \mathbf{L})(\mathbf{X} \mathbf{w} - \hat{\mathbf{y}})$ . Since this is a sparse matrix multiplication, the complexity of this step is  
 573  $O(m + nd)$ . Putting things together, we have the time complexity of calling inner problem solver is  
 574  $\tilde{O}\left((m + nd)\lambda^2 \log \frac{1}{\beta}\right)$ .

575 From Lemma 5.1, the number of outer iterations needed is  $\tilde{O}(\kappa_o \log \frac{1}{\epsilon})$ , where  $\kappa_o$  is the condition  
576 number of the outer problem. Moreover, Lemma 5.1 also indicates that we require  $\beta \leq \frac{\epsilon^{1/2}}{25\kappa(\mathbf{X})\lambda}$ , so  
577 solving inner problem takes  $\tilde{O}(\lambda^2 \log 1/\beta) = \tilde{O}(\lambda^2 \log 1/\epsilon)$ . In each outer iteration, we need to call  
578 inner problem solver constant times, and as well as performing constant matrix vector multiplications  
579 whose complexity is  $O(nd)$ . Therefore the overall running time of solving the training problem is  
580  $\tilde{O}\left[\kappa_o(\lambda^2(m+nd)+nd)\left(\log \frac{1}{\epsilon}\right)^2\right]$ .

## 581 C Proof of Theoretical Results

582 We first note that since  $\hat{\mathbf{L}}$  is normalized, all the eigenvalues of  $\hat{\mathbf{L}}$  are in the range  $[0, 2)$ . Therefore,  
583 we have  $\sigma_{\min}(\mathbf{I} + \lambda\hat{\mathbf{L}}) = 1$  and  $\sigma_{\max}(\mathbf{I} + \lambda\hat{\mathbf{L}}) \leq 1 + 2\lambda$ . In the whole paper we view  $\lambda$  as a large  
584 value (i.e.  $\lambda \gg 1$ ). In practice, it is possible to use a small  $\lambda$ , in which case the algorithm works no  
585 worse than the case where  $\lambda = 1$ . Therefore the  $\lambda$  used in the paper should actually be understood as  
586  $\max\{\lambda, 1\}$ . With this assumption, below we assume  $\sigma_{\max}(\mathbf{I} + \lambda\hat{\mathbf{L}}) \leq 3\lambda$  for convenience.

### 587 C.1 Descent Lemma

588 In this subsection, we introduce a descent lemma that we will use to analyze the convergence rate.  
589 This is a standard result in convex optimization, and we refer interested readers to [3] for more details.

590 **Lemma C.1.** *If  $\ell : \mathbb{R}^d \rightarrow \mathbb{R}$  is  $L$ -Lipschitz smooth and  $c$ -strongly convex, and we have a sequence of  
591 points  $\{\mathbf{w}^{(t)}\}_{t=1}^T$  in  $\mathbb{R}^d$  such that*

$$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} - \eta \mathbf{g}^{(t)}, \quad (\text{C.1})$$

592 where  $\|\mathbf{g}^{(t)} - \nabla \ell(\mathbf{w}^{(t)})\|_2 \leq \gamma \|\nabla \ell(\mathbf{w}^{(t)})\|_2$  and  $\gamma < 1$ ,  $\eta = \frac{1-\gamma}{(1+\gamma)^2 L}$ , then we have

$$\ell(\mathbf{w}^{(T)}) - \ell^* \leq \left[1 - \kappa^{-1} \left(\frac{1-\gamma}{1+\gamma}\right)^2\right]^T \left[\ell(\mathbf{w}^{(0)}) - \ell^*\right], \quad (\text{C.2})$$

593 where  $\ell^* = \inf_{\mathbf{w} \in \mathbb{R}^d} \ell(\mathbf{w})$  and  $\kappa = \frac{L}{c}$ .

594 *Proof.* From the condition  $\|\mathbf{g}^{(t)} - \nabla \ell(\mathbf{w}^{(t)})\|_2 \leq \gamma \|\nabla \ell(\mathbf{w}^{(t)})\|_2$ , we have that

$$(1-\gamma)\|\mathbf{g}^{(t)}\|_2 \leq \|\nabla \ell(\mathbf{w}^{(t)})\|_2 \leq (1+\gamma)\|\mathbf{g}^{(t)}\|_2 \quad (\text{C.3})$$

595 and

$$-2\langle \mathbf{g}^{(t)}, \nabla \ell(\mathbf{w}^{(t)}) \rangle = \|\mathbf{g}^{(t)} - \nabla \ell(\mathbf{w}^{(t)})\|_2^2 - \|\mathbf{g}^{(t)}\|_2^2 - \|\nabla \ell(\mathbf{w}^{(t)})\|_2^2 \quad (\text{C.4})$$

$$\leq [(\gamma^2 - 1) - (1 - \gamma)^2] \|\nabla \ell(\mathbf{w}^{(t)})\|_2^2 \quad (\text{C.5})$$

$$= 2(\gamma - 1) \|\nabla \ell(\mathbf{w}^{(t)})\|_2^2 \quad (\text{C.6})$$

596 From Lipschitz smoothness, we have

$$\ell(\mathbf{w}^{(t+1)}) - \ell(\mathbf{w}^{(t)}) \leq -\eta \langle \mathbf{g}^{(t)}, \nabla \ell(\mathbf{w}^{(t)}) \rangle + \frac{1}{2} L \eta^2 \|\mathbf{g}^{(t)}\|_2^2 \quad (\text{C.7})$$

$$\leq \left(\eta(\gamma - 1) + \frac{1}{2} L \eta^2 (1 + \gamma)^2\right) \|\nabla \ell(\mathbf{w}^{(t)})\|_2^2. \quad (\text{C.8})$$

597 It's not hard to show that the optimal  $\eta$  for eq.(C.8) is  $\eta = \frac{1-\gamma}{L(1+\gamma)^2}$ . Substituting  $\eta = \frac{1-\gamma}{L(1+\gamma)^2}$  to  
598 eq.(C.8) and use convexity we can get

$$\ell(\mathbf{w}^{(t+1)}) - \ell^* \leq \left[\ell(\mathbf{w}^{(t)}) - \ell^*\right] - \frac{(1-\gamma)^2}{2L(1+\gamma)^2} \|\nabla \ell(\mathbf{w}^{(t)})\|_2^2 \quad (\text{C.9})$$

$$\leq \left[\ell(\mathbf{w}^{(t)}) - \ell^*\right] - \frac{c(1-\gamma)^2}{L(1+\gamma)^2} \left[\ell(\mathbf{w}^{(t)}) - \ell^*\right] \quad (\text{C.10})$$

599 By induction we have

$$\ell(\mathbf{w}^{(T)}) - \ell^* \leq \left[1 - \kappa^{-1} \left(\frac{1-\gamma}{1+\gamma}\right)^2\right]^T [\ell(\mathbf{w}^{(0)}) - \ell^*]. \quad (\text{C.11})$$

600

□

## 601 C.2 Bound of Loss Value by Gradient

602 In this subsection, we derive an inequality that allows us to bound the value of loss function by the  
603 norm of gradient.

604 **Lemma C.2.** *If  $\ell : \mathbb{R}^d \rightarrow \mathbb{R}$  is a  $\mathcal{C}$ -strongly convex and smooth function, and  $\mathbf{w}^* =$   
605  $\arg \min_{\mathbf{w} \in \mathbb{R}^d} \ell(\mathbf{w})$  is a global optimal, then for any  $\mathbf{w} \in \mathbb{R}^d$  and  $\epsilon \in (0, 1)$ , we have*

$$\ell(\mathbf{w}) \leq \max \left\{ (1 + \epsilon)\ell(\mathbf{w}^*), \frac{1}{\epsilon\mathcal{C}} \|\nabla \ell(\mathbf{w})\|^2 \right\}. \quad (\text{C.12})$$

606 *Proof.* Using the inequality (4.12) from [3], we have

$$\ell(\mathbf{w}) \leq \frac{1}{2\mathcal{C}} \|\nabla \ell(\mathbf{w})\|^2 + \ell(\mathbf{w}^*). \quad (\text{C.13})$$

607 If  $\ell(\mathbf{w}) \geq (1 + \epsilon)\ell(\mathbf{w}^*)$ , we have

$$\ell(\mathbf{w}) \leq \frac{1}{2\mathcal{C}} \|\nabla \ell(\mathbf{w})\|^2 + \frac{1}{1 + \epsilon} \ell(\mathbf{w}). \quad (\text{C.14})$$

608 Shifting the terms in eq.(C.14) gives  $\ell(\mathbf{w}) \leq \frac{1+\epsilon}{2\mathcal{C}\epsilon} \|\nabla \ell(\mathbf{w})\|^2 \leq \frac{1}{\epsilon\mathcal{C}} \|\nabla \ell(\mathbf{w})\|^2$ , which proves the  
609 claim.

610

□

## 611 C.3 Proof of Lemma 5.1

612 Given a point  $\mathbf{w} \in \mathbb{R}^d$ , we first consider the error between the estimated gradient  $\widetilde{\nabla \ell_i}(\mathbf{w})$  and the  
613 true gradient  $\nabla \ell_i(\mathbf{w})$ . Below we denote  $\mathbf{H} = \mathbf{I} + \lambda \hat{\mathbf{L}}$  and  $\mathbf{z} = \mathbf{X}\mathbf{w} - \hat{\mathbf{y}}_i$ . It's not hard to notice that  
614  $\ell_i(\mathbf{w}) = \frac{1}{2} \|\mathbf{H}^{-1}\mathbf{z}\|^2$  and  $\nabla \ell_i(\mathbf{w}) = \mathbf{X}^\top \mathbf{H}^{-2}\mathbf{z}$ . Moreover, we denote  $\ell_i^* = \inf_{\mathbf{w}} \ell_i^*$  as the optimal  
615 value of  $\ell_i$ .

616 we have

$$\left\| \widetilde{\nabla \ell_i}(\mathbf{w}) - \nabla \ell_i(\mathbf{w}) \right\| = \left\| \mathbf{X}^\top [\mathcal{S}(\mathbf{H}^{-1}\mathbf{z}) - \mathbf{H}^{-2}\mathbf{z} + \mathcal{S}(\mathcal{S}(\mathbf{z})) - \mathcal{S}(\mathbf{H}^{-1}\mathbf{z})] \right\| \quad (\text{C.15})$$

$$\leq \left\| \mathbf{X}^\top [\mathcal{S}(\mathbf{H}^{-1}\mathbf{z}) - \mathbf{H}^{-2}\mathbf{z}] \right\| + \left\| \mathbf{X}^\top [\mathcal{S}(\mathcal{S}(\mathbf{z})) - \mathbf{H}^{-1}\mathbf{z}] \right\|. \quad (\text{C.16})$$

617 As shown, the gradient error can be decomposed into two terms, namely  $\left\| \widetilde{\nabla \ell_i}(\mathbf{w}) - \nabla \ell_i(\mathbf{w}) \right\| =$   
618  $E_1 + E_2$ , where  $E_1 = \left\| \mathbf{X}^\top [\mathcal{S}(\mathbf{H}^{-1}\mathbf{z}) - \mathbf{H}^{-2}\mathbf{z}] \right\|$  and  $E_2 = \left\| \mathbf{X}^\top [\mathcal{S}(\mathcal{S}(\mathbf{z})) - \mathbf{H}^{-1}\mathbf{z}] \right\|$ .  
619 Below we analyze each part separately.

620 Notice that, as we assumed  $\mathbf{X}$  is full-rank,  $\ell_i$  is a strongly convex function with parameter  $\mathcal{C} =$   
621  $\sigma_{\min}(\mathbf{X}^\top \mathbf{H}^{-2} \mathbf{X}) \geq \frac{\sigma_{\min}(\mathbf{X})^2}{9\lambda^2}$ . Let  $\mathcal{D} = \{\mathbf{w} \in \mathbb{R}^d \mid \ell_i(\mathbf{w}) \leq (1 + \epsilon)\ell_i^*\}$ . If  $\mathbf{w} \in \mathcal{D}$ , then  $\mathbf{w}$  is  
622 already a good enough solution. Below we assume  $\mathbf{w} \notin \mathcal{D}$ , in which case by Lemma C.2, we have  
623  $\ell_i(\mathbf{w}) \leq \frac{1}{\epsilon\mathcal{C}} \|\nabla \ell_i(\mathbf{w})\|^2$ , which means

$$\|\mathbf{H}^{-1}\mathbf{z}\|^2 \leq \frac{2}{\epsilon\mathcal{C}} \|\mathbf{X}^{-1}\mathbf{H}^{-2}\mathbf{z}\|^2 \leq \frac{18\lambda^2}{\sigma_{\min}(\mathbf{X})^2\epsilon}. \quad (\text{C.17})$$

624 For  $E_1$ , we have

$$E_1 = \|\mathbf{X}^\top [\mathcal{S}(\mathbf{H}^{-1}\mathbf{z}) - \mathbf{H}^{-2}\mathbf{z}]\|_2 \quad (\text{C.18})$$

$$\leq \sigma_{\max}(\mathbf{X}) \|\mathcal{S}(\mathbf{H}^{-1}\mathbf{z}) - \mathbf{H}^{-2}\mathbf{z}\|_2 \quad (\text{C.19})$$

$$\leq \sigma_{\max}(\mathbf{X}) \|\mathcal{S}(\mathbf{H}^{-1}\mathbf{z}) - \mathbf{H}^{-2}\mathbf{z}\|_{\mathbf{H}} \quad (\text{C.20})$$

$$\leq \sigma_{\max}(\mathbf{X})\mu \|\mathbf{H}^{-2}\mathbf{z}\|_{\mathbf{H}} \quad (\text{C.21})$$

$$\leq \sigma_{\max}(\mathbf{X})\mu \|\mathbf{H}^{-1}\mathbf{z}\|_2 \quad (\text{C.22})$$

$$\leq \sqrt{\frac{18}{\epsilon}} \kappa(\mathbf{X})\lambda\mu \|\mathbf{X}^\top \mathbf{H}^{-2}\mathbf{z}\|_2 \quad (\text{C.23})$$

$$\leq 5\epsilon^{-1/2} \kappa(\mathbf{X})\lambda\mu \|\nabla \ell_i(\mathbf{w})\|_2. \quad (\text{C.24})$$

625 For  $E_2$ , we have

$$E_2 = \|\mathbf{X}^\top [\mathcal{S}(\mathcal{S}(\mathbf{z}) - \mathbf{H}^{-1}\mathbf{z})]\|_2 \quad (\text{C.25})$$

$$\leq \sigma_{\max}(\mathbf{X})(1+\mu) \|\mathbf{H}^{-1}(\mathcal{S}(\mathbf{z}) - \mathbf{H}^{-1}\mathbf{z})\|_{\mathbf{H}} \quad (\text{C.26})$$

$$\leq \sigma_{\max}(\mathbf{X})(1+\mu) \|\mathcal{S}(\mathbf{z}) - \mathbf{H}^{-1}\mathbf{z}\|_{\mathbf{H}} \quad (\text{C.27})$$

$$\leq \sigma_{\max}(\mathbf{X})(1+\mu)\mu \|\mathbf{H}^{-1}\mathbf{z}\|_{\mathbf{H}} \quad (\text{C.28})$$

$$\leq \sqrt{\frac{18}{\epsilon}} \kappa(\mathbf{X})(1+\mu)\mu\lambda\sqrt{3}\lambda \|\mathbf{X}^\top \mathbf{H}^{-2}\mathbf{z}\|_2 \quad (\text{C.29})$$

$$\leq 20\epsilon^{-1/2} \kappa(\mathbf{X})\mu\lambda^2 \|\nabla \ell_i(\mathbf{w}_i)\|_2. \quad (\text{C.30})$$

626 Combine the bound of  $E_1$  and  $E_2$ , we have

$$\|\widetilde{\nabla \ell_i(\mathbf{w}_i)} - \nabla \ell_i(\mathbf{w}_i)\| \leq E_1 + E_2 \leq 25\epsilon^{-1/2} \kappa(\mathbf{X})\mu\lambda^2 \|\nabla \ell_i(\mathbf{w}_i)\|_2. \quad (\text{C.31})$$

627 Now, consider the optimization. Let  $\{\mathbf{w}_i^{(t)}\}_{t=1}^T$  be a sequence such that

$$\mathbf{w}_i^{(t+1)} = \mathbf{w}_i^{(t)} - \eta \widetilde{\nabla \ell_i(\mathbf{w}_i^{(t)})}. \quad (\text{C.32})$$

628 Let  $\kappa$  be the condition number of  $\mathbf{X}^\top \mathbf{H}^{-2} \mathbf{X}$  and let  $\ell_i^*$  be the optimal value of  $\ell_i$ . Let  $\gamma =$   
 629  $25\epsilon^{-1/2} \kappa(\mathbf{X})\lambda^2\mu \leq \frac{1}{2}$ . From Lemma C.1, we have with a proper value of  $\eta$ ,

$$\ell_i(\mathbf{w}_i^{(T)}) - \ell_i^* \leq \left(1 - \left(\frac{1-\gamma}{1+\gamma}\right)^2 \kappa^{-1}\right)^T [\ell_i(\mathbf{w}_i^{(0)}) - \ell_i^*] \quad (\text{C.33})$$

$$\leq (1 - (9\kappa)^{-1})^T [\ell_i(\mathbf{w}_i^{(0)}) - \ell_i^*]. \quad (\text{C.34})$$

630 Therefore, the number of iterations  $T$  required to achieve  $\epsilon$  error rate is

$$T = O\left(\left(\log \frac{1}{1 - (9\kappa)^{-1}}\right)^{-1} \log \frac{1}{\epsilon}\right) = O\left(\kappa \log \frac{1}{\epsilon}\right). \quad (\text{C.35})$$

#### 631 C.4 Proof of Lemma 5.2

632 The basic idea of Algorithm 2 is to use ridge leverage score sampling methods to obtain a spectral  
 633 sparsifier. Let  $\hat{\mathbf{B}} = \mathbf{B}\mathbf{D}^{-1/2}$  be the normalized incidence matrix and  $\hat{\mathbf{b}}_i$  the  $i$ -th row of  $\hat{\mathbf{B}}$ . Given  
 634  $\lambda^{-1} > 0$ , for  $i \in \{1, 2, \dots, m\}$ , the  $i$ -th ridge leverage score is defined as

$$l_i := \hat{\mathbf{b}}_i^\top (\hat{\mathbf{L}} + \lambda^{-1} \mathbf{I})^{-1} \hat{\mathbf{b}}_i \quad (\text{C.36})$$

$$= \hat{\mathbf{b}}_i^\top (\hat{\mathbf{L}} + \lambda^{-1} \mathbf{I})^{-1} \hat{\mathbf{L}} (\hat{\mathbf{L}} + \lambda^{-1} \mathbf{I})^{-1} \hat{\mathbf{b}}_i + \lambda^{-1} \hat{\mathbf{b}}_i^\top (\hat{\mathbf{L}} + \lambda^{-1} \mathbf{I})^{-2} \hat{\mathbf{b}}_i \quad (\text{C.37})$$

$$= \|\hat{\mathbf{B}} (\hat{\mathbf{L}} + \lambda^{-1} \mathbf{I})^{-1} \hat{\mathbf{b}}_i\|^2 + \lambda^{-1} \|(\hat{\mathbf{L}} + \lambda^{-1} \mathbf{I})^{-1} \hat{\mathbf{b}}_i\|^2. \quad (\text{C.38})$$

635 It is not affordable to compute all the  $m$  ridge leverage scores exactly. Therefore, we first use  
636 Johnson–Lindenstrauss lemma to reduce the dimension. Recall that in Algorithm 2 we define  
637  $\mathbf{\Pi}_1 \in \mathbb{R}^{k \times m}$  and  $\mathbf{\Pi}_2 \in \mathbb{R}^{k \times n}$  to be Gaussian sketches (that is, each entry of the matrices are i.i.d  
638 Gaussian random variables  $\mathcal{N}(0, 1/k)$ ). We set  $k = O(\log m)$ , and by using Lemma C.3 we have  
639 the following claim holds with probability at least  $1 - \frac{1}{8n}$ :

$$\|\mathbf{\Pi}_1 \hat{\mathbf{B}}(\hat{\mathbf{L}} + \lambda^{-1} \mathbf{I})^{-1} \hat{\mathbf{b}}_i\| \approx_{2^{1/4-1}} \|\hat{\mathbf{B}}(\hat{\mathbf{L}} + \lambda^{-1} \mathbf{I})^{-1} \hat{\mathbf{b}}_i\| \quad \text{for all } i \in \{1, 2, \dots, m\}. \quad (\text{C.39})$$

640 Similarly, the following claim also holds with probability at least  $1 - \frac{1}{8n}$ :

$$\|\mathbf{\Pi}_2(\hat{\mathbf{L}} + \lambda^{-1} \mathbf{I})^{-1} \hat{\mathbf{b}}_i\| \approx_{2^{1/4-1}} \|(\hat{\mathbf{L}} + \lambda^{-1} \mathbf{I})^{-1} \hat{\mathbf{b}}_i\| \quad \text{for all } i \in \{1, 2, \dots, m\}. \quad (\text{C.40})$$

641 By summing over above two inequalities (after squared) and taking an union bound, with probability  
642  $1 - \frac{1}{4n}$  we have

$$\|\mathbf{\Pi}_1 \hat{\mathbf{B}}(\hat{\mathbf{L}} + \lambda^{-1} \mathbf{I})^{-1} \hat{\mathbf{b}}_i\|^2 + \lambda^{-1} \|\mathbf{\Pi}_2(\hat{\mathbf{L}} + \lambda^{-1} \mathbf{I})^{-1} \hat{\mathbf{b}}_i\|^2 \approx_{\sqrt{2-1}} l_i \quad \text{for all } i \in \{1, 2, \dots, m\}. \quad (\text{C.41})$$

643 However it is still too expensive to compute  $\mathbf{\Pi}_1 \hat{\mathbf{B}}(\hat{\mathbf{L}} + \lambda^{-1} \mathbf{I})^{-1}$  and  $\mathbf{\Pi}_2(\hat{\mathbf{L}} + \lambda^{-1} \mathbf{I})^{-1}$ , since  
644 computing  $(\hat{\mathbf{L}} + \lambda^{-1} \mathbf{I})^{-1}$  itself takes prohibitive  $O(n^3)$  time. Instead, notice that  $\hat{\mathbf{L}} + \lambda^{-1} \mathbf{I}$  is a  
645 SDD matrix, thus we can apply the SDD solver to the  $k$  columns of matrix  $(\mathbf{\Pi}_1 \hat{\mathbf{B}})^\top$  and  $(\mathbf{\Pi}_2)^\top$   
646 respectively. According to Lemma A.1, there is a linear operator  $\text{SSolve}_\epsilon(\hat{\mathbf{L}} + \lambda^{-1} \mathbf{I}; \mathbf{x})$  that runs in  
647 time  $\tilde{O}(\text{nnz}(\hat{\mathbf{L}} + \lambda^{-1} \mathbf{I}) \cdot \log 1/\epsilon) = \tilde{O}(m \log 1/\epsilon)$  such that for any  $\mathbf{x}^\top \in \mathbb{R}^n$ , it outputs  $\tilde{\mathbf{x}}$  that satisfies  
648  $\|\tilde{\mathbf{x}} - \mathbf{x}^\top (\hat{\mathbf{L}} + \lambda^{-1} \mathbf{I})^{-1}\|_{\hat{\mathbf{L}} + \lambda^{-1} \mathbf{I}} \leq \epsilon \|\mathbf{x}^\top (\hat{\mathbf{L}} + \lambda^{-1} \mathbf{I})\|_{\hat{\mathbf{L}} + \lambda^{-1} \mathbf{I}}$ . For our purpose we set  $\epsilon = 2^{1/4}$ .  
649 Denote  $\mathbf{B}_S$  as the matrix obtained by applying each column of  $(\mathbf{\Pi}_1 \hat{\mathbf{B}})^\top$  to  $\text{SSolve}_\epsilon(\hat{\mathbf{L}} + \lambda^{-1} \mathbf{I}; \mathbf{x})$ ,  
650 that is, the  $j$ -th row of  $\text{SSolve}_{2^{1/4}}(\hat{\mathbf{L}} + \lambda^{-1} \mathbf{I}; (\mathbf{\Pi}_1 \hat{\mathbf{B}})_j^\top)$ . Similarly we denote  $\mathbf{\Pi}_S$  as the matrix with  
651  $j$ -th row equals to  $\text{SSolve}_{2^{1/4}}(\hat{\mathbf{L}} + \lambda^{-1} \mathbf{I}; (\mathbf{\Pi}_2)_j^\top)$ . By using Lemma A.1 to both solvers we have

$$\|\mathbf{B}_S \hat{\mathbf{b}}_i\|^2 + \lambda^{-1} \|\mathbf{\Pi}_S \hat{\mathbf{b}}_i\|^2 \approx_{\sqrt{2-1}} \|\mathbf{\Pi}_1 \hat{\mathbf{B}}(\hat{\mathbf{L}} + \lambda^{-1} \mathbf{I})^{-1} \hat{\mathbf{b}}_i\|^2 + \lambda^{-1} \|\mathbf{\Pi}_2(\hat{\mathbf{L}} + \lambda^{-1} \mathbf{I})^{-1} \hat{\mathbf{b}}_i\|^2 \quad (\text{C.42})$$

652 for all  $i \in \{1, 2, \dots, m\}$ . By eq.(C.41) and eq.(C.42), if we set  $\tilde{l}_i := \|\mathbf{B}_S \hat{\mathbf{b}}_i\|^2 + \lambda^{-1} \|\mathbf{\Pi}_S \hat{\mathbf{b}}_i\|^2$ ,  
653 then with probability  $1 - \frac{1}{4n}$ , we obtain all the approximation of ridge leverage scores  $\{\tilde{l}_i\}_{i=1}^m$  such  
654 that  $\tilde{l}_i \approx_{1/2} l_i$  holds for all  $i$ . Notice that since  $\mathbf{B}_S, \mathbf{\Pi}_S \in \mathbb{R}^{k \times n}$ , and that each  $\hat{\mathbf{b}}_i$  only contains 2  
655 non-zero entries, thus it takes  $\tilde{O}(k \cdot m)$  to pre-compute  $\mathbf{B}_S$  and  $\mathbf{\Pi}_S$ , and takes  $O(2k \cdot m)$  to compute  
656  $\{\mathbf{B}_S \hat{\mathbf{b}}_i, \mathbf{\Pi}_S \hat{\mathbf{b}}_i\}_{i=1}^m$ . To summarize, computing all  $\tilde{l}_i$  takes  $\tilde{O}(km) = \tilde{O}(m \log m) = \tilde{O}(m)$ . With  
657 these ridge leverage score approximations, we apply Lemma C.4 to matrix  $\hat{\mathbf{B}}$  with choice  $\delta = 1/4n$ .  
658 By setting  $\tilde{\mathbf{L}} := \hat{\mathbf{B}}^\top \mathbf{S}^\top \mathbf{S} \hat{\mathbf{B}}$  we have  $\tilde{\mathbf{L}} + \lambda^{-1} \mathbf{I} \approx_\epsilon \hat{\mathbf{L}} + \lambda^{-1} \mathbf{I}$  holds with probability  $1 - 1/4n$ . By  
659 applying another union bound we obtain our final result:

$$\tilde{\mathbf{L}} + \lambda^{-1} \mathbf{I} \approx_\epsilon \hat{\mathbf{L}} + \lambda^{-1} \mathbf{I} \quad \text{with probability } 1 - \frac{1}{2n}. \quad (\text{C.43})$$

660 Finally, according to Lemma C.4, the number of edges of  $\tilde{\mathbf{L}}$  is  $s = Cn_\lambda \log(n^2)/\epsilon^2 =$   
661  $O(n_\lambda \log n/\epsilon^2)$ . Since the last step of computing  $\tilde{\mathbf{L}} = (\mathbf{S} \hat{\mathbf{B}})^\top (\mathbf{S} \hat{\mathbf{B}})$  only takes  $O(s) = \tilde{O}(n_\lambda/\epsilon^2)$   
662 due to the sparsity of  $\hat{\mathbf{B}}$ , the overall time complexity of Algorithm 2 is  $\tilde{O}(m + n_\lambda/\epsilon^2)$ .

663 **Lemma C.3** (Johnson–Lindenstrauss, [12]). *Let  $\mathbf{\Pi} \in \mathbb{R}^{k \times n}$  be Gaussian sketch matrix with each  
664 entry independent and equal to  $\mathcal{N}(0, 1/k)$  where  $\mathcal{N}(0, 1)$  denotes a standard Gaussian random  
665 variable. If we choose  $k = O\left(\frac{\log(1/\delta)}{\epsilon^2}\right)$ , then for any vector  $\mathbf{x} \in \mathbb{R}^n$ , with probability  $1 - \delta$  we have*

$$(1 - \epsilon) \|\mathbf{x}\| \leq \|\mathbf{\Pi} \mathbf{x}\| \leq (1 + \epsilon) \|\mathbf{x}\|. \quad (\text{C.44})$$

666 **Lemma C.4** (Spectral approximation, [11]). *Let  $\mathbf{S}$  be an  $s \times m$  subsampling matrix with probabilities  
667  $p_i = \tilde{\ell}_i/Z$  where  $\tilde{\ell}_i \approx_{1/2} \ell_i$  and  $Z$  is the normalization constant. If we have  $s \geq Cn_\lambda \log(n/\delta)/\epsilon^2$   
668 for some constant  $C > 0$  and  $\epsilon, \delta \in (0, 1/2]$ , then we have*

$$\hat{\mathbf{B}}^\top \mathbf{S}^\top \mathbf{S} \hat{\mathbf{B}} + \lambda^{-1} \mathbf{I} \approx_\epsilon \hat{\mathbf{B}}^\top \hat{\mathbf{B}} + \lambda^{-1} \mathbf{I} \quad (\text{C.45})$$

669 holds with probability  $1 - \delta$ . Here  $n_\lambda = \text{Tr}[\hat{\mathbf{L}}(\hat{\mathbf{L}} + \lambda^{-1} \mathbf{I})^{-1}]$ .

670 **C.5 Proof of Lemma 5.3**

671 We first prove a lemma which is related to the approximation rate of squared matrices.

672 **Lemma C.5.** *Suppose that  $\Sigma$  and  $\tilde{\Sigma}$  are two  $n \times n$  PD matrices, and  $\tilde{\Sigma} \approx_{\frac{\beta}{\kappa}} \Sigma$ , where  $\kappa$  is the*  
 673 *condition number of  $\Sigma$  and  $\beta \in (0, \frac{1}{8})$ , then we have*

$$\tilde{\Sigma}^2 \approx_{8\beta} \Sigma^2. \quad (\text{C.46})$$

674 *Proof.* Let  $\epsilon = \frac{\beta}{\kappa}$ . The condition  $\Sigma \approx_{\epsilon} \tilde{\Sigma}$  implies that  $\|\Sigma - \tilde{\Sigma}\| \leq \epsilon \|\Sigma\|$ . We have

$$\|(\Sigma - \tilde{\Sigma})\mathbf{x}\| = \|(\Sigma - \tilde{\Sigma})\Sigma^{-1}\Sigma\mathbf{x}\| \quad (\text{C.47})$$

$$\leq \|(\Sigma - \tilde{\Sigma})\Sigma^{-1}\| \times \|\Sigma\mathbf{x}\| \quad (\text{C.48})$$

$$\leq \|(\Sigma - \tilde{\Sigma})\| \|\Sigma^{-1}\| \times \|\Sigma\mathbf{x}\| \quad (\text{C.49})$$

$$\leq \epsilon \|\Sigma\| \|\Sigma^{-1}\| \times \|\Sigma\mathbf{x}\| \quad (\text{C.50})$$

$$= \epsilon\kappa \|\Sigma\mathbf{x}\| \quad (\text{C.51})$$

$$= \beta \|\Sigma\mathbf{x}\|. \quad (\text{C.52})$$

675 For any  $\mathbf{x} \in \mathbb{R}^n$ , by triangle inequality we have

$$\|\Sigma\mathbf{x}\| - \|(\Sigma - \tilde{\Sigma})\mathbf{x}\| \leq \|\tilde{\Sigma}\mathbf{x}\| \leq \|\Sigma\mathbf{x}\| + \|(\Sigma - \tilde{\Sigma})\mathbf{x}\|. \quad (\text{C.53})$$

676 Subtracting the inequality  $\|(\Sigma - \tilde{\Sigma})\mathbf{x}\| \leq \beta \|\Sigma\mathbf{x}\|$  we derived before into eq.(C.53) and squaring  
 677 all sides, we have

$$(1 - \beta)^2 \mathbf{x}^\top \Sigma^2 \mathbf{x} \leq \mathbf{x}^\top \tilde{\Sigma}^2 \mathbf{x} \leq (1 + \beta)^2 \mathbf{x}^\top \Sigma^2 \mathbf{x}. \quad (\text{C.54})$$

678 Since  $0 < \beta < \frac{1}{8}$ , the claim is proved.  $\square$

679 Let  $\mathbf{T} = \mathbf{X}^\top (\mathbf{I} + \lambda \hat{\mathbf{L}})^{-2} \mathbf{X}$  be the true Hessian. Let  $\tilde{\mathbf{T}} = \mathbf{X}^\top (\mathbf{I} + \lambda \tilde{\mathbf{L}})^{-2} \mathbf{X}$  be the approxi-  
 680 mated Hessian using sparsified Laplacian  $\tilde{\mathbf{L}}$ . By Lemma 5.2,  $\mathbf{I} + \lambda \tilde{\mathbf{L}} \approx_{\frac{\beta}{3\lambda}} \mathbf{I} + \lambda \mathbf{L}$  with probability  
 681 at least  $1 - \frac{1}{2n}$ . From Lemma C.5, we have  $\tilde{\mathbf{T}} \approx_{8\beta} \mathbf{T}$ .

682 Next, we show that  $\mathbf{Q}^\top \mathbf{Q} \approx_{O(\beta)} \mathbf{T}$ . Let  $\tilde{\mathbf{H}} = \mathbf{I} + \lambda \tilde{\mathbf{L}}$ . Let  $\mathcal{S}$  be the operator defined by  
 683  $\text{SSolve}_{\frac{\beta}{\sqrt{3\lambda}}}(\tilde{\mathbf{H}}, \cdot)$ , i.e.  $\mathbf{q}_j = \mathcal{S}(\mathbf{x}_j)$  where  $\mathbf{q}_j$  and  $\mathbf{x}_j$  are the  $j$ -th column of  $\mathbf{Q}$  and  $\mathbf{X}$  respectively.  
 684 From Lemma A.1 we have for any  $\mathbf{z} \in \mathbb{R}^d$ ,

$$\|\mathbf{Q}\mathbf{z} - \tilde{\mathbf{H}}^{-1}\mathbf{X}\mathbf{z}\|_2^2 = \left\| \sum_{j=1}^d z_j (\mathcal{S}(\mathbf{x}_j) - \tilde{\mathbf{H}}\mathbf{x}_j) \right\|_2^2 \quad (\text{C.55})$$

$$= \left\| \mathcal{S} \left( \sum_{j=1}^d z_j \mathbf{x}_j \right) - \tilde{\mathbf{H}} \left( \sum_{j=1}^d z_j \mathbf{x}_j \right) \right\|_2^2 \quad (\text{C.56})$$

$$\leq \frac{\beta^2}{3\lambda} \|\tilde{\mathbf{H}}^{-1}\mathbf{X}\mathbf{z}\|_{\tilde{\mathbf{H}}}^2 \quad (\text{C.57})$$

$$\leq \beta^2 \|\tilde{\mathbf{H}}^{-1}\mathbf{X}\mathbf{z}\|_2^2. \quad (\text{C.58})$$

685 Therefore we have

$$(1 - \beta) \|\tilde{\mathbf{H}}^{-1}\mathbf{X}\mathbf{z}\|_2 \leq \|\mathbf{Q}\mathbf{z}\| \leq (1 + \beta) \|\tilde{\mathbf{H}}^{-1}\mathbf{X}\mathbf{z}\|_2, \quad (\text{C.59})$$

686 and this is equivalent to

$$(1 - \beta)^2 \mathbf{z}^\top (\mathbf{X} \tilde{\mathbf{H}}^{-2} \mathbf{X}) \mathbf{z} \leq \mathbf{z}^\top \mathbf{Q}^\top \mathbf{Q} \mathbf{z} \leq (1 + \beta)^2 \mathbf{z}^\top (\mathbf{X}^\top \tilde{\mathbf{H}}^{-2} \mathbf{X}) \mathbf{z}. \quad (\text{C.60})$$

687 Notice that  $\mathbf{X}^\top \tilde{\mathbf{H}}^{-2} \mathbf{X} = \tilde{\mathbf{T}}$ . We conclude that  $\mathbf{Q}^\top \mathbf{Q} \approx_{2\beta + \beta^2} \tilde{\mathbf{T}}$ . When  $\beta < \frac{1}{8}$ , we have  
688  $\mathbf{Q}^\top \mathbf{Q} \approx_{4\beta} \tilde{\mathbf{T}}$ .

689 Lastly, from Lemma A.2 and the discussions in Appendix A.3, we have there exists a constant  
690  $C' \in (0, 1/4)$ , such that

$$\mathbf{P} = \tilde{\mathbf{Q}}^\top \tilde{\mathbf{Q}} \approx_{C'} \mathbf{Q}^\top \mathbf{Q} \quad (\text{C.61})$$

691 with probability at least  $1 - \frac{1}{2n}$ .

692 Put the results above together and we get

$$\mathbf{P} \approx_{12\beta + C'} \mathbf{T}. \quad (\text{C.62})$$

693 Notice that  $12\beta + C' < \frac{1}{2}$ . Using a union bound can prove that the fail probability of this whole  
694 process is bounded by  $\frac{1}{n}$ .

## 695 C.6 Proof of Lemma 5.4

696 Let  $\mathbf{F} = (\mathbf{I} + \lambda \hat{\mathbf{L}}) \mathbf{X}$ . We have the Hessian of  $\ell'$  is

$$\nabla^2 \ell' = \mathbf{P}^{-\frac{1}{2}} \mathbf{F}^\top \mathbf{F} \mathbf{P}^{-\frac{1}{2}}. \quad (\text{C.63})$$

697 From the condition that  $\mathbf{P} \approx_{c_0} \mathbf{F}^\top \mathbf{F}$ , we have  $\mathbf{F}^\top \mathbf{F} \preceq (1 + c_0) \mathbf{P}$ , which implies

$$\mathbf{P}^{-1/2} \mathbf{F}^\top \mathbf{F} \mathbf{P}^{-1/2} \preceq (1 + c_0) \mathbf{I}. \quad (\text{C.64})$$

698 Similarly, since  $\mathbf{P} \preceq (1 + c_0) \mathbf{F}^\top \mathbf{F}$ , we have

$$(\mathbf{F}^\top \mathbf{F})^{-1} \preceq (1 + c_0) \mathbf{P}^{-1}, \quad (\text{C.65})$$

699 which implies

$$\left[ \lambda_{\min}(\mathbf{P}^{-1/2} \mathbf{F}^\top \mathbf{F} \mathbf{P}^{-1/2}) \right]^{-1} = \lambda_{\max} \left( \mathbf{P}^{1/2} (\mathbf{F}^\top \mathbf{F})^{-1} \mathbf{P}^{1/2} \right) \quad (\text{C.66})$$

$$\leq 1 + c_0. \quad (\text{C.67})$$

700 .

701 Notice that  $\ell'$  and  $\ell$  have the same global optimal value, let it be  $\ell^*$ . For any  $\mathbf{w}'$  satisfies  $\ell'(\mathbf{w}') - \ell^* \leq$   
702  $\gamma$ , we have

$$\ell \left( \mathbf{P}^{1/2} \mathbf{w}' \right) - \ell^* = \ell'(\mathbf{w}') - \ell^* \leq \gamma. \quad (\text{C.68})$$

703 Therefore, if  $\mathbf{w}'$  is a solution of  $\ell'$  with  $\epsilon$  error rate, then  $\mathbf{P}^{1/2} \mathbf{w}'$  is a solution of  $\ell$  with  $\epsilon$  error rate.

## 704 C.7 Proof of Theorem 5.1

705 As we noted in the main paper, Algorithm 1 is composed by two components: constructing the  
706 preconditioner  $\mathbf{P}$  and applying it to the optimization.

707 **Constructing the Preconditioner.** From Lemma 5.2, the first step that applies spectral sparsifier  
708 to get  $\tilde{\mathbf{L}}$  requires  $\tilde{O}(m)$  time and the number of non-zero entries in  $\tilde{\mathbf{L}}$  is  $\tilde{O}(n_\lambda \lambda^2)$ . By Lemma A.1,  
709 running the SDD solver with error rate  $\frac{\beta}{3\lambda}$  in the second step requires  $\tilde{O}(n_\lambda \lambda^2)$  time. Notice that  
710 since  $\mathbf{X}$  is an  $n \times d$  matrix, we actually need to run SDD solver for  $d$  times, and this introduces  
711 another  $d$  factor in the time complexity.

712 As we noted in Appendix A.3, applying the Hadamard transformation to  $\mathbf{RQ}$  requires  $\tilde{O}(d)$  time,  
713 and applying the subsampling requires  $\tilde{O}(n)$  time. Since  $\tilde{\mathbf{Q}} \in \mathbb{R}^{s \times d}$ , calculating  $\mathbf{P} = \tilde{\mathbf{Q}}^\top \mathbf{Q}$   
714 requires  $O(d^2 s)$  time. Since  $s = \tilde{O}(d)$ , this step takes  $\tilde{O}(d^3)$  time. We use brute force to calculate  
715  $\mathbf{P}' = \mathbf{P}^{-1/2}$ , and this takes  $O(d^3)$  time.

716 As a summary, constructing the preconditioner  $\mathbf{P}'$  takes  $\tilde{O}(n_\lambda \lambda^2 d + nd + d^3)$  time. Notice that we  
717 only perform this step once during the whole algorithm.

718 **Performing the Iterations.** Next we consider the time required for each iteration. We calculate  
719  $\mathbf{u}^{(t)}$  from right to left and it takes  $O(nd)$  time. Next, we need to perform two SDD solvers with error  
720 rate  $\mu$ . From Lemma A.1, it takes  $\tilde{O}\left[m \log\left(\frac{1}{\mu}\right)\right]$ , which is  $\tilde{O}\left[m \log\left(\frac{1}{\epsilon}\right)\right]$ . Calculating  $\mathbf{g}^{(t)}$  and  
721  $\mathbf{w}^{(t)}$  is straight-forward and takes  $\tilde{O}(nd)$  time.

722 To conclude, performing each iteration requires  $\tilde{O}\left[m \log\left(\frac{1}{\epsilon}\right) + nd\right]$  time. By Lemma 5.1 and  
723 Lemma 5.3, with a proper step size, the number of iterations needed for solving outer problem is  
724  $\tilde{O}(\log 1/\epsilon)$ .

725 Combining the analysis above together, to overall complexity is

$$\tilde{O}\left[n_\lambda \lambda^2 d + nd + d^3 + \left(m \log\left(\frac{1}{\epsilon}\right) + nd\right) \log\left(\frac{1}{\epsilon}\right)\right] = \tilde{O}\left(n_\lambda \lambda^2 d + d^3 + (m + nd) (\log 1/\epsilon)^2\right), \quad (\text{C.69})$$

726 which proves the claim.

## 727 D Further Discussions

728 In this section, we extend some of the discussions in the main paper.

### 729 D.1 An Analysis of CE Loss v.s. MSE Loss

730 Although HERTA is derived from MSE loss, the experiment result shows it also works on CE loss. In  
731 this subsection we provide an analysis showing the similarity of the gradient and Hessian of TWIRLS  
732 on CE loss and MSE loss, to offer a intuitive explanation why a method that is derived from MSE  
733 loss can work on CE loss.

734 In this section, for a node  $u \in \{1, 2, \dots, n\}$ , we use  $\mathbf{y}^{(u)} \in \mathbb{R}^c$  to represent the  $u$ -th row of  $\mathbf{Y}$   
735 (notice we use super-script here to distinguish from  $\mathbf{y}_i$  used before),  $\mathbf{h}_u \in \mathbb{R}^n$  to represent the  $u$ -th  
736 row of  $(\mathbf{I} + \lambda \hat{\mathbf{L}})^{-1}$ . For  $p \in \{1, 2, \dots, d\}$ , we use  $\mathbf{x}_p \in \mathbb{R}^n$  to denote the  $p$ -th column of  $\mathbf{X}$ . For  
737  $i \in \{1, 2, \dots, c\}$ , we use  $y_i^{(u)}$  to denote the  $i$ -th entry of  $\mathbf{y}^{(u)}$  (or in other words the  $(u, i)$ -th entry of  
738  $\mathbf{Y}$ ), and  $w_{p,i}$  to denote the  $(p, i)$ -th entry of  $\mathbf{W}$ .

739 The MSE loss of TWIRLS can be decomposed into summation of sub-losses of each node, i.e.

$$\ell(\mathbf{W}) = \sum_{u=1}^n \frac{1}{2} \left\| \mathbf{h}_u^\top \mathbf{X} \mathbf{W} - \mathbf{y}^{(u)} \right\|_{\mathcal{F}}^2 = \sum_{u=1}^n \ell^{(u)}(\mathbf{W}), \quad (\text{D.1})$$

740 where  $\ell^{(u)}(\mathbf{W}) = \frac{1}{2} \left\| \mathbf{W}^\top \mathbf{X}^\top \mathbf{h}_u - \mathbf{y}^{(u)} \right\|_{\mathcal{F}}^2$ . For a specific class  $i \in \{1, 2, \dots, c\}$ , we have the  
741 gradient of  $\ell^{(u)}$  w.r.t.  $\mathbf{w}_i$  is

$$\frac{\partial \ell^{(u)}}{\partial \mathbf{w}_i} = \mathbf{X}^\top \mathbf{h}_u \mathbf{h}_u^\top \mathbf{X} \mathbf{w}_i - \mathbf{X}^\top \mathbf{h}_u \times y_i^{(u)}. \quad (\text{D.2})$$

742 It is not hard to see from eq.(D.2) that the Hessian of  $\ell^{(u)}$  with respect to  $\mathbf{w}_i$  is

$$\nabla_{\mathbf{w}_i}^2 \ell^{(u)}(\mathbf{W}) = \mathbf{X}^\top \mathbf{h}_u \mathbf{h}_u^\top \mathbf{X}. \quad (\text{D.3})$$

743 For classification tasks, the target  $\mathbf{y}^{(u)}$ -s are one-hot vectors representing the class of the node. Notice  
744 when we calculate cross entropy loss, we use softmax to normalize it before feeding it into the loss  
745 function. In the following for a vector  $\mathbf{v} \in \mathbb{R}^c$  whose  $i$ -th entry is  $v_i$ , we define

$$\text{softmax}(\mathbf{v})_i = \frac{\exp(v_i)}{\sum_{j=1}^c \exp(v_j)}. \quad (\text{D.4})$$

746 Suppose the  $u$ -th node belongs to class  $k$ , then the cross entropy loss of the  $u$ -th node is defined as

$$\text{CE}^{(u)}(\mathbf{W}) = -\log \text{softmax}(\mathbf{h}_u^\top \mathbf{X} \mathbf{W})_k \quad (\text{D.5})$$

$$= -\mathbf{h}_u^\top \mathbf{X} \mathbf{w}_k + \log \sum_{j=1}^c \exp(\mathbf{h}_u^\top \mathbf{X} \mathbf{w}_j), \quad (\text{D.6})$$

747 The gradient of  $\text{CE}^{(u)}$  w.r.t.  $\mathbf{w}_i$  is

$$\frac{\partial \text{CE}^{(u)}}{\partial \mathbf{w}_i} = -\delta_{i,k} \times \mathbf{X}^\top \mathbf{h}_u + \frac{\mathbf{X}^\top \mathbf{h}_u \exp(\mathbf{h}_u^\top \mathbf{X} \mathbf{w}_i)}{\sum_{j=1}^c \exp(\mathbf{h}_u^\top \mathbf{X} \mathbf{w}_j)} \quad (\text{D.7})$$

$$= \mathbf{X}^\top \mathbf{h}_u \text{softmax}(\mathbf{h}_u^\top \mathbf{X} \mathbf{W})_i - \mathbf{X}^\top \mathbf{h}_u \times y_i^{(u)}. \quad (\text{D.8})$$

748 By taking another partial differentiation, we have

$$\frac{\partial^2 \text{CE}^{(u)}}{\partial w_{p,i} \partial w_{q,i}} = \frac{(\mathbf{x}_p^\top \mathbf{h}_u) (\mathbf{x}_q^\top \mathbf{h}_u) \exp(\mathbf{h}_u^\top \mathbf{X} \mathbf{w}_i)}{\sum_{j=1}^c \exp(\mathbf{h}_u^\top \mathbf{X} \mathbf{w}_j)} - \frac{(\mathbf{x}_p^\top \mathbf{h}_u) (\mathbf{x}_q^\top \mathbf{h}_u) [\exp(\mathbf{h}_u^\top \mathbf{X} \mathbf{w}_i)]^2}{\left[ \sum_{j=1}^c \exp(\mathbf{h}_u^\top \mathbf{X} \mathbf{w}_j) \right]^2} \quad (\text{D.9})$$

$$= (\mathbf{x}_p^\top \mathbf{h}_u) (\mathbf{x}_q^\top \mathbf{h}_u) [s_i - s_i^2], \quad (\text{D.10})$$

749 where  $s_i = \text{softmax}(\mathbf{h}_u^\top \mathbf{X} \mathbf{W})_i$ . Rewriting eq.(D.10) into matrix form, we have

$$\nabla_{\mathbf{w}_i}^2 \text{CE}^{(u)}(\mathbf{W}) = (s_i - s_i^2) \mathbf{X}^\top \mathbf{h}_u \mathbf{h}_u^\top \mathbf{X}. \quad (\text{D.11})$$

750 Comparing eq.(D.2) and eq.(D.8), we can notice that the gradient of MSE loss and CE loss are  
 751 essentially the same except the term  $\mathbf{h}_u^\top \mathbf{X} \mathbf{w}_i$  in eq.(D.2) is normalized by softmax in eq.(D.8).  
 752 Additionally, by comparing the Hessian of MSE loss eq.(D.3) and the Hessian of CE loss eq.(D.11),  
 753 we have the latter is a rescaled version of the former. These similarities intuitively explains why our  
 754 method can be also used on CE loss.

## 755 D.2 The Unavoidable $\lambda^2$ in the Running Time

756 Lemma C.5 essentially states the following idea: even when two PD matrices are very similar (in  
 757 terms of spectral approximation), if they are ill-conditioned, their square can be very different. In  
 758 this subsection, we prove that the bound obtained in Lemma C.5 is strict up to constant factors by  
 759 explicitly constructing a worst case. Notice that, in this subsection we discuss the issue of squaring  
 760 spectral approximators in a broader context, thus in this subsection we possibly overload some  
 761 symbols used before to simplify the notation.

762 First we consider another definition of approximation rate: for two PD matrices  $\Sigma$  and  $\tilde{\Sigma}$ , we define  
 763 the approximation rate as

$$\psi(\Sigma, \tilde{\Sigma}) = \max \left\{ \left\| \Sigma^{-1/2} \tilde{\Sigma} \Sigma^{-1/2} \right\|, \left\| \tilde{\Sigma}^{-1/2} \Sigma \tilde{\Sigma}^{-1/2} \right\| \right\}. \quad (\text{D.12})$$

764 This definition is easier to calculate in the scenario considered in this subsection, and can be easily  
 765 translated to the definition we used in Section 3: when  $\epsilon = \psi(\Sigma, \tilde{\Sigma}) \in (0, 1)$ , we have  $\tilde{\Sigma} \approx_\epsilon \Sigma$ ,  
 766 and when  $\psi(\Sigma, \tilde{\Sigma})$  is larger than 1 we don't have an approximation in the form defined in Section 3.

767 Let  $\gamma > 1$  and  $\delta \in (0, 1)$  be real numbers. Define

$$\Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} \gamma & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, \quad (\text{D.13})$$

768 and

$$\tilde{\Sigma} = \begin{bmatrix} \sqrt{1-\delta^2} & -\delta \\ \delta & \sqrt{1-\delta^2} \end{bmatrix} \begin{bmatrix} \gamma & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} \sqrt{1-\delta^2} & \delta \\ -\delta & \sqrt{1-\delta^2} \end{bmatrix}. \quad (\text{D.14})$$

769 It's not hard to compute the error rate  $\psi(\Sigma, \tilde{\Sigma}) = \left\| \Sigma^{-1/2} \tilde{\Sigma} \Sigma^{-1/2} \right\| \approx \Theta(\gamma \delta^2)$ . While if we  
 770 consider the squared matrices, we have  $\psi(\Sigma^2, \tilde{\Sigma}^2) = \left\| \Sigma^{-1} \tilde{\Sigma}^2 \Sigma^{-1} \right\| \approx \gamma^2 \delta^2$ . In this case,  
 771  $\psi(\Sigma^2, \tilde{\Sigma}^2)$  is larger than  $\psi(\Sigma, \tilde{\Sigma})$  by a factor of  $\gamma$ , which is the condition number of  $\Sigma$ . When  $\gamma$   
 772 is very large and  $\delta$  is very small, we can have a small  $\psi(\Sigma, \tilde{\Sigma})$  while large  $\psi(\Sigma^2, \tilde{\Sigma}^2)$ , which .

773 **A More General Construction.** For a PD matrix  $\Sigma$  and its spectral approximation  $\tilde{\Sigma}$ , we call  
774  $\frac{\psi(\Sigma^2, \tilde{\Sigma}^2)}{\psi(\Sigma, \tilde{\Sigma})}$  the Squared Error Rate. As noted above, in the worst case the squared error rate can be as  
775 large as the condition number of  $\Sigma$ . However, the construction above is limited to  $2 \times 2$  matrices.  
776 Now we construct a more general worst case of the squared error rate and perform a loose analysis.  
777 Although not rigorously proved, the construction and the analysis suggest the origination of large  
778 squared error rates: it approaches the upper bound (condition number) when the eigenspace of the  
779 two matrices are very well aligned but not exactly the same.

780 Let  $\mathbf{A}$  be an ill-conditioned matrix with all eigenvalues very large except one eigenvalue equals to  
781 1, and the eigenvalues of  $\mathbf{B}$  are all closed to the eigenvalues of  $\mathbf{A}$ . Specifically, Let the SVD of  
782  $\mathbf{A}$  be  $\mathbf{A} = \mathbf{U} \text{diag}(\boldsymbol{\lambda}) \mathbf{U}^\top$ , where  $\boldsymbol{\lambda} = [\lambda_1 \ \lambda_2 \ \cdots \ \lambda_n]$ , and suppose  $\lambda_n = 1$  and  $\lambda_k \gg 1, \forall k \leq$   
783  $n - 1$ . For simplicity we just let  $\mathbf{A}$  and  $\mathbf{B}$  have the same eigenvalues. Suppose the SVD of  $\mathbf{B}$  is  
784  $\mathbf{B} = \mathbf{V} \text{diag}(\boldsymbol{\lambda}) \mathbf{V}^\top$  and let  $\boldsymbol{\Lambda} = \text{diag}(\boldsymbol{\lambda})$ . We denote the one sided error rate of  $\mathbf{A}$  and  $\mathbf{B}$  by  $\epsilon$ , i.e.  
785  $\epsilon = \|\mathbf{A}^{-1/2} \mathbf{B} \mathbf{A}^{-1/2}\|$ . We have

$$\epsilon = \left\| \boldsymbol{\Lambda}^{-1/2} \mathbf{U}^\top \mathbf{V} \boldsymbol{\Lambda} \mathbf{V}^\top \mathbf{U} \boldsymbol{\Lambda}^{-1/2} \right\| \quad (\text{D.15})$$

$$= \left\| \boldsymbol{\Lambda}^{-1/2} \mathbf{W} \boldsymbol{\Lambda} \mathbf{W}^\top \boldsymbol{\Lambda}^{-1/2} \right\|, \quad (\text{D.16})$$

786 where  $\mathbf{W} = \mathbf{U}^\top \mathbf{V}$ .

787 Since  $\lambda_n = 1$  and for all  $k \leq n - 1, \lambda_k \gg 1$ , we have

$$\boldsymbol{\Lambda}^{-\frac{1}{2}} = \begin{bmatrix} \lambda_1^{-1/2} & & & \\ & \lambda_2^{-1/2} & & \\ & & \ddots & \\ & & & \lambda_n^{-1/2} \end{bmatrix} \approx \begin{bmatrix} 0 & & & \\ & 0 & & \\ & & \ddots & \\ & & & 1 \end{bmatrix}, \quad (\text{D.17})$$

788 and therefore

$$\epsilon = \left\| \boldsymbol{\Lambda}^{-1/2} \mathbf{W} \boldsymbol{\Lambda} \mathbf{W}^\top \boldsymbol{\Lambda}^{-1/2} \right\| \approx (\mathbf{W} \boldsymbol{\Lambda} \mathbf{W}^\top)_{n,n} = \sum_{k=1}^n \mathbf{W}_{k,n}^2 \lambda_k, \quad (\text{D.18})$$

789 where  $\mathbf{A}_{i,j}$  represents the  $(i, j)$ -th entry of a matrix  $\mathbf{A}$ .

790 For simplicity, we assume all  $\lambda_k (k \leq n - 1)$  are approximate equal, i.e.  $\lambda_1 \approx \lambda_2 \cdots \approx \lambda_{n-1} = \gamma$ .  
791 Then we have

$$\epsilon = \sum_{k=1}^n \mathbf{W}_{k,n}^2 \lambda_k \quad (\text{D.19})$$

$$\approx \gamma \left[ \sum_{k=1}^{n-1} \mathbf{W}_{k,n}^2 + \frac{1}{\gamma} \mathbf{W}_{n,n}^2 \right] \quad (\text{D.20})$$

$$\approx \gamma \sum_{k=1}^{n-1} \mathbf{W}_{k,n}^2 \quad (\text{D.21})$$

$$= \gamma (1 - \mathbf{W}_{n,n}^2). \quad (\text{D.22})$$

792 Recall  $\mathbf{W} = \mathbf{U}^\top \mathbf{V}$ , we have  $\mathbf{W}_{n,n} = \mathbf{u}_n^\top \mathbf{v}_n$ , where  $\mathbf{u}_n$  and  $\mathbf{v}_n$  are the  $n$ -th column of  $\mathbf{U}$  and  $\mathbf{V}$   
793 respectively. Thus we have

$$\epsilon \approx \gamma \left[ 1 - (\mathbf{u}_n^\top \mathbf{v}_n)^2 \right]. \quad (\text{D.23})$$

794 Now we can see the spectral approximation error  $\epsilon$  is determined by two terms: the condition number  
795  $\gamma$  and the matchness of the eigenvectors corresponds to small eigenvalues, which is evaluated by  
796  $\left[ 1 - (\mathbf{u}_n^\top \mathbf{v}_n)^2 \right]$ . **When  $\gamma$  is very large but  $\left[ 1 - (\mathbf{u}_n^\top \mathbf{v}_n)^2 \right]$  is small,  $\mathbf{B}$  can still be a spectral  
797 **approximation of  $\mathbf{A}$ .** For example if  $(1 - \mathbf{u}_n^\top \mathbf{v}_n) \approx \gamma^{-1}$ , we can get a spectral approximation  
798 error  $\epsilon \approx 1$ . However, after we square the matrices, the eigenvalues will also get squared, but  
799 eigenvalues remains unchanged. That will enlarge the spectral approximation error by a factor of  $\gamma$ ,  
800 i.e.  $\|\mathbf{A}^{-1} \mathbf{B}^2 \mathbf{A}^{-1}\| \approx \gamma^2 \gamma^{-1} = \gamma$ , which becomes very large. In this case the squared error rate is  
801  $\gamma$ , the condition number of  $\mathbf{A}$ .**

### 802 D.3 Applying Graph Sparsification in Each Iteration

803 As mentioned in the main paper, unlike most existing work, in our algorithm we don't sparsify the  
 804 graph in each training iteration. The proof of Lemma 5.1 (see Appendix C.3) suggests the reason  
 805 why performing graph sparsification in each iteration can lead to suboptimal running time. In this  
 806 subsection we illustrate this claim in detail. The intuition is that, in order to ensure convergence of  
 807 the training, we require a small error rate in the gradient estimation, which is of the order  $\epsilon^{1/2}$  as we  
 808 have showed in Appendix C.3. It is acceptable for the SDD solver because the running time of the  
 809 SDD solver only logarithmly depends on the error rate. However, if we sparsify the graph, to obtain  
 810 an  $\epsilon^{1/2}$  error rate we will need to sample  $O(n_\lambda/\epsilon)$  edges, which grows linearly with  $1/\epsilon$ , and can be  
 811 large especially when  $\epsilon$  is very small. Below is a more detailed analysis.

812 Consider in the for-loop of Algorithm 1 we replace the  $\hat{\mathbf{L}}$  by a sparsified version  $\mathbf{L}' = \text{Sparsify}_\omega(\hat{\mathbf{L}})$ ,  
 813 where  $\omega \in (0, 1)$ . Let  $\mathbf{H}' = \mathbf{I} + \lambda\mathbf{L}'$ . From Lemma 5.2, we have  $\mathbf{H}' \approx_\omega \mathbf{H}$ . Consider  $E_1$  defined  
 814 in Appendix C.3, it now becomes

$$E'_1 = \|\mathbf{X}^\top [\mathcal{S}(\mathbf{H}'^{-1}\mathbf{z}) - \mathbf{H}^{-2}\mathbf{z}]\| \quad (\text{D.24})$$

$$\leq \sigma_{\max}(\mathbf{X}) \|\mathcal{S}(\mathbf{H}'^{-1}\mathbf{z}) - \mathbf{H}^{-2}\mathbf{z}\|. \quad (\text{D.25})$$

815 Here we can not proceed by using the fact that  $\mathcal{S}$  is a linear solver, because now  $\mathcal{S}$  is not a linear  
 816 solver for  $\mathbf{H}$  but for  $\mathbf{H}'$ , thus we will have to split the error term again:

$$E'_1 \leq \sigma_{\max}(\mathbf{X}) \|\mathcal{S}(\mathbf{H}'^{-1}\mathbf{z}) - \mathbf{H}'^{-2}\mathbf{z}\| + \sigma_{\max}(\mathbf{X}) \|\mathbf{H}^{-2}\mathbf{z} - \mathbf{H}'^{-2}\mathbf{z}\|. \quad (\text{D.26})$$

817 Of the two terms on the right-hand side of eq.(D.26), the first one can be bounded with a similar  
 818 method used in Appendix C.3, and the second term is bounded by

$$\sigma_{\max}(\mathbf{X}) \|\mathbf{H}^{-2}\mathbf{z} - \mathbf{H}'^{-2}\mathbf{z}\| \leq \sigma_{\max}(\mathbf{X}) \|\mathbf{I} - \mathbf{H}^{-1}\mathbf{H}'^{-2}\mathbf{H}^{-1}\| \|\mathbf{H}^{-2}\mathbf{z}\| \quad (\text{D.27})$$

$$\leq \sigma_{\max}(\mathbf{X}) \omega \|\mathbf{H}^{-2}\mathbf{z}\| \quad (\text{D.28})$$

$$\leq \sqrt{\frac{8}{\epsilon}} \kappa(\mathbf{X}) \lambda \omega \|\nabla \ell_i(\mathbf{w})\|. \quad (\text{D.29})$$

819 Therefore, in order to obtain  $E'_1 \leq \|\nabla \ell_i(\mathbf{w})\|$ , we at least require  $\omega \leq \frac{\epsilon^{-1/2}}{\sqrt{8\kappa(\mathbf{X})\lambda}} = O(\epsilon^{-1/2})$ . From  
 820 Lemma 5.2, the number of edges in the sparsified graph  $\omega$  error rate is  $O(n_\lambda/\epsilon)$ . A similar analysis  
 821 can be also applied to  $E_2$ , and by repeating the proof in Appendix C.3, we obtain an overall running  
 822 time bound

$$O(m + n_\lambda \lambda^2 d + d^3 + n_\lambda \epsilon^{-1} \log 1/\epsilon), \quad (\text{D.30})$$

823 which, although eliminates the  $m(\log \frac{1}{\epsilon})^2$  term, introduces an extra  $n_\lambda \epsilon^{-1} \log \frac{1}{\epsilon}$  term, and is usually  
 824 worse than the original bound we derived in Theorem 5.1, especially when requiring a relatively  
 825 small  $\epsilon$ .

826 That being said, although not very likely, when  $\epsilon$  and  $l$  or  $\lambda$  are large, it is possible that directly  
 827 sparsifying the graph in each iteration is beneficial. Considering this, we can adopt a mixed strategy:  
 828 when  $n_\lambda \epsilon^{-1} \leq m \log(1/\epsilon)$ , we sparsify the graph in each iteration, otherwise we don't. This leads to  
 829 the following overall running time bound:

$$O\left[m + n_\lambda \lambda^2 d + d^3 + \min\left\{m \log \frac{1}{\epsilon}, n_\lambda \epsilon^{-1}\right\} \log \frac{1}{\epsilon}\right], \quad (\text{D.31})$$

830 which is slightly better than the one we presented in Theorem 5.1.

### 831 D.4 Possible Directions for Extending HERTA to More Complex Models

832 The assumption that  $f(\mathbf{X}; \mathbf{W}) = \mathbf{X}\mathbf{W}$  provides us with a convenience that the Hessian of this  
 833 model is a constant matrix. Therefore in Algorithm 1 we only need to calculate the preconditioner  $\mathbf{P}$   
 834 for one time. However, if  $f(\mathbf{X}; \mathbf{W})$  is implemented by a non-linear network, then the Hessian will  
 835 change by the time and might be hard to calculate, which will be a key challenge to using a more  
 836 complex  $f$ . We note that this can be possibly addressed by constructing a linear approximation of  
 837  $f$  using its Jacobian. In each iteration, we can use the Jacobian to replace the  $\mathbf{X}$  used in current

838 version of HERTA, and recalculate  $\mathbf{P}$  at each iteration. Since the convergence is fast, we only need  
 839 to recalculate the Jacobian for a small number of iterations, so should not bring massive change to  
 840 the running time of the algorithm.

841 The attention mechanism of TWIRLS in [46] is achieved by adding a concave penalty function to  
 842 each summand of the the  $\text{Tr}(\mathbf{Z}^\top \hat{\mathbf{L}} \mathbf{Z})$  term in eq.(4.1). For specific penalty functions, it might  
 843 still possible to find a inner problem solver, as long as the problem stay convex. We note that this  
 844 depends on concrete implementation of the penalty term used. Investigating how to fast solve the  
 845 inner problem under various penalty functions should also be an important problem for future study.

## 846 E Implementation Details

847 In the experiments, the datasets are loaded and processed using the DGL package [40]. We use the  
 848 original inner problem solver in [46] since it is not computation bottleneck. It can be in principle  
 849 replaced by any implementation of SDD solvers.

850 For the calculating the gradient of the preconditioned model, we presented a calculation method in  
 851 Algorithm 1 which maintains the lowest computational complexity. In preliminary experiments, we  
 852 tested this calculation method with using the autograd module in pytorch<sup>3</sup> and verified that they have  
 853 the same output, and similar computational efficiency on real world datasets (again in practice this  
 854 step is not a bottleneck). Therefore, we simply use pytorch autograd module to compute gradients in  
 855 experiments.

856 Using the pytorch autograd module also enables us to apply HERTA on various loss functions and  
 857 optimizers: we only need to perform the preconditioning and indicate the loss function. The gradient  
 858 and optimization algorithm will be automatically realized by pytorch.

859 In order to ensure a fair comparison and prevent confounding factors, we don't using any common  
 860 training regularization techniques like weight decay or dropout. For each setting, we repeat the  
 861 experiment with learning rates in  $\{0.001, 0.01, 0.1, 1, 10\}$  and choose the trial which the training loss  
 862 does not explode and with the lowest final training loss to report in the main paper.

## 863 F Additional Experiment Results

864 In this section, we present additional experiment results.

### 865 F.1 Experiments with Larger $\lambda$

866 All the experiments presented in the main paper are with  $\lambda = 1$ . In this subsection, we present results  
 867 with  $\lambda = 20$ . See Figures 3 and 4 for results with MSE loss and CE loss respectively. The results  
 868 supports our observation in the main paper that HERTA works consistently well on all settings.

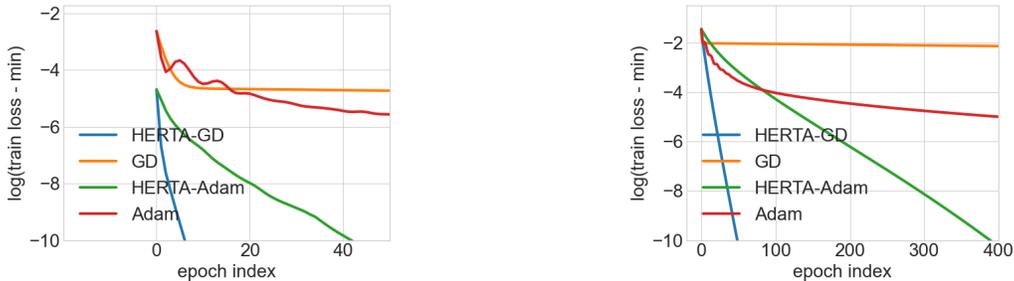


Figure 3: The training loss comparison between HERTA and standard optimizers on MSE loss with  $\lambda = 20$ . Dataset used from left to right: ogbn-arxiv, pubmed.

<sup>3</sup><https://pytorch.org/>

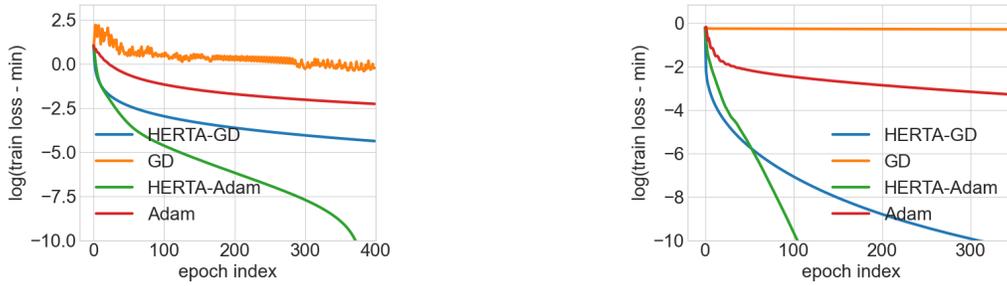


Figure 4: The training loss comparison between HERTA and standard optimizers on CE loss with  $\lambda = 20$ . Dataset used from left to right: ogbn-arxiv, pubmed.

869 **F2 Additional Results on Cora and Citeseer**

870 In this subsection we present results on Cora and Citeseer, which are other citation datasets used in  
 871 [46]. See Figure 5 and Figure 7 for results with  $\lambda = 1$  and  $\lambda = 20$  on Cora respectively. See Figure 6  
 872 and Figure 8 for results with  $\lambda = 1$  and  $\lambda = 20$  on Citeseer respectively. It is clear that the results on  
 873 Cora is consistent with our observation on other datasets.

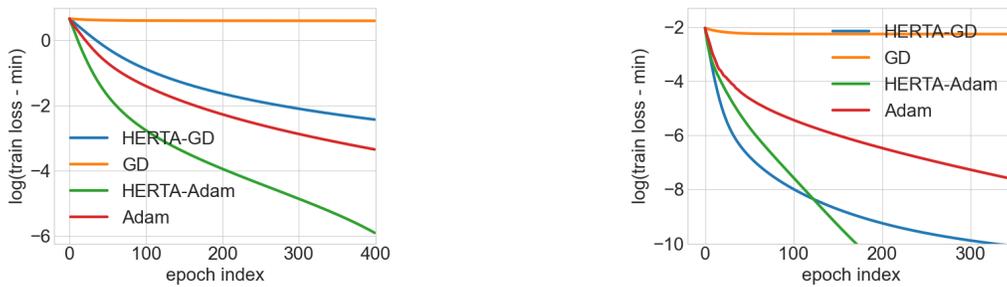


Figure 5: The training loss comparison between HERTA and standard optimizers on Cora with  $\lambda = 1$ . Left: CE loss. Right: MSE loss.

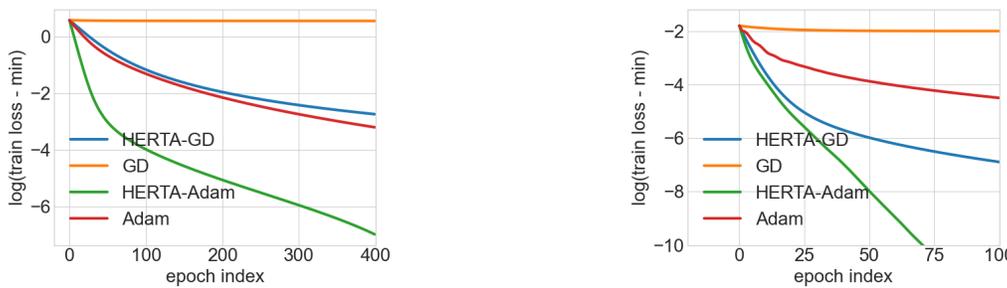


Figure 6: The training loss comparison between HERTA and standard optimizers on Cora with  $\lambda = 1$ . Left: CE loss. Right: MSE loss.

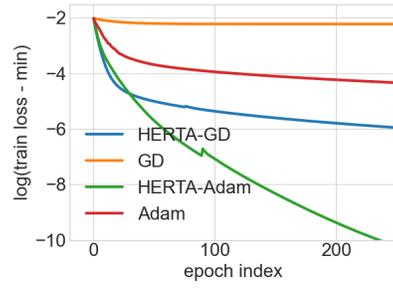
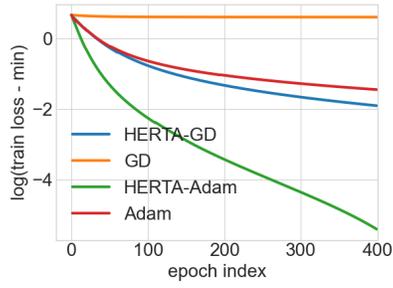


Figure 7: The training loss comparison between HERTA and standard optimizers on Cora with  $\lambda = 20$ . Left: CE loss. Right: MSE loss.

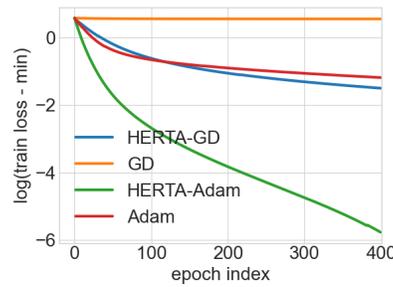
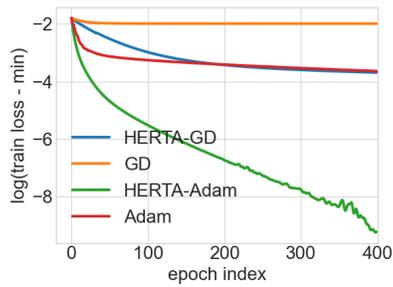


Figure 8: The training loss comparison between HERTA and standard optimizers on Citeseer with  $\lambda = 20$ . Left: CE loss. Right: MSE loss.