

Rethinking Smoothness in Node Features Learned by Graph Convolutional Networks

Anonymous authors
Paper under double-blind review

Abstract

The pioneering works of Oono and Suzuki (ICLR 2020) and Cai and Wang (arXiv:2006.13318) initiated the analysis of feature smoothness in graph convolutional networks (GCNs), uncovering a strong empirical connection between node classification accuracy and the ratio of smooth to non-smooth feature components. However, it remains unclear how to effectively control this ratio in learned node features to enhance classification performance. Furthermore, deep GCNs with ReLU or leaky ReLU activations tend to suppress non-smooth feature components. In this paper, we introduce a novel strategy to enable GCNs to learn node features with **controllable smoothness**, thereby improving node classification accuracy. Our method comprises three core components: (1) deriving a geometric relationship between the inputs and outputs of ReLU and leaky ReLU activations; (2) augmenting the standard message-passing mechanism in graph convolutional layers with a learnable term for efficient smoothness modulation; and (3) theoretically analyzing the attainable smooth-to-non-smooth ratios under the proposed augmented propagation. Extensive experiments demonstrate that our approach substantially enhances node classification performance across GCNs and related architectures.

1 Introduction

Let $G = (V, E)$ be an undirected graph with node set $V = \{v_i\}_{i=1}^n$ and edge set E . Let $\mathbf{A} \in \mathbb{R}^{n \times n}$ be the adjacency matrix of the graph G , where $A_{ij} = \mathbf{1}_{(i,j) \in E}$ with $\mathbf{1}$ being the indicator function. Moreover, let \mathbf{G} be the augmented normalized adjacency matrix:

$$\mathbf{G} := (\mathbf{D} + \mathbf{I})^{-\frac{1}{2}}(\mathbf{I} + \mathbf{A})(\mathbf{D} + \mathbf{I})^{-\frac{1}{2}} = \tilde{\mathbf{D}}^{-\frac{1}{2}}\tilde{\mathbf{A}}\tilde{\mathbf{D}}^{-\frac{1}{2}}, \quad (1)$$

where $\mathbf{I} \in \mathbb{R}^{n \times n}$ is the identity matrix, $\mathbf{D} \in \mathbb{R}^{n \times n}$ is a diagonal matrix with $D_{ii} = \sum_{j=1}^n A_{ij}$. Starting from the initial node features $\mathbf{H}^0 := [(\mathbf{h}_1^0)^\top, \dots, (\mathbf{h}_n^0)^\top]^\top \in \mathbb{R}^{d \times n}$, the graph convolutional network (GCN) (Kipf & Welling, 2017) updates node features via a cascade of the following graph convolutional layer (GCL):

$$\mathbf{H}^l = \sigma(\mathbf{W}^l \mathbf{H}^{l-1} \mathbf{G}), \quad (2)$$

where σ is the activation function (typically ReLU), and $\mathbf{W}^l \in \mathbb{R}^{d \times d}$ is a learnable weight matrix. GCL smooths node features by aggregating information from neighbors, a property shown to benefit node classification (cf. (Li et al., 2018; Wu et al., 2019; Chen et al., 2020a)), echoing the idea of energy-based semi-supervised learning (cf. (Zhu et al., 2003; Zhou et al., 2003)). However, accurate node classification requires a balance between smooth and non-smooth feature components (Oono & Suzuki, 2020). Beyond GCNs, many other graph neural networks (GNNs) have been proposed using diverse mechanisms, e.g., spectral methods (Defferrard et al., 2016), spatial methods (Gilmer et al., 2017; Veličković et al., 2018), sampling methods (Hamilton et al., 2017a; Ying et al., 2018), and the attention mechanism (Veličković et al., 2018). Many other GNNs can be found in e.g., (Hamilton et al., 2017b; Battaglia et al., 2018; Wu et al., 2020; Zhou et al., 2020; Hamilton, 2020).

It has been observed that deep GCNs perform significantly worse than shallow models (Chen et al., 2020a). The node features of deep GCNs tend to be identical over each connected component of the graph; this

phenomenon is referred to as **over-smoothing** (Li et al., 2018; Nt & Maehara, 2019; Oono & Suzuki, 2020; Cai & Wang, 2020; Chen et al., 2020a; Wu et al., 2023b), which not only occurs for GCN but also for many other GNNs (Hamilton et al., 2017a; Gilmer et al., 2017; Wu et al., 2023a). Intuitively, each GCL smooths neighboring node features; stacking many smoothing layers will inevitably homogenize node features. Algorithms have been developed to alleviate over-smoothing, including decoupling prediction and message passing (Gasteiger et al., 2019), skip connection and batch normalization (Kawamoto et al., 2018; Chen et al., 2019; 2020b; Scholkemper et al., 2024), graph sparsification (Rong et al., 2020), jumping knowledge (Xu et al., 2018), PairNorm (Zhao & Akoglu, 2020), implicit layers (Chamberlain et al., 2021; Thorpe et al., 2022; Gu et al., 2020; Baker et al., 2023; 2024), and controlling the Dirichlet energy of node features (Zhou et al., 2021).

From a theoretical standpoint, Oono & Suzuki (2020) show that, as the depth of a GCN with ReLU activation increases, node features collapse to the eigenspace \mathcal{M} associated with the largest eigenvalue 1 of the matrix \mathbf{G} in equation 1. They also empirically demonstrate that node classification accuracy is intricately linked to the ratio between smooth and non-smooth components of node features, i.e., their projections onto \mathcal{M} and its orthogonal complement \mathcal{M}^\perp , respectively. Their findings suggest that both components are essential for accurate node classification, though the optimal ratio remains unknown and task-dependent. Complementing this, Cai & Wang (2020) prove that the Dirichlet energy—a measure of feature smoothness—also vanishes as GCN depth increases.

A key insight from the analyses in (Oono & Suzuki, 2020; Cai & Wang, 2020) is that both ReLU and leaky ReLU activations reduce the distance of node features to the eigenspace \mathcal{M} and lower their Dirichlet energy. However, Cai & Wang (2020) point out that this notion of over-smoothing—based solely on proximity to \mathcal{M} or low Dirichlet energy—is inappropriate. Cai & Wang (2020) propose that smoothness should instead be assessed via a **normalized smoothness**, such as Dirichlet energy normalized by the feature magnitude. This perspective aligns with the ratio of smooth to non-smooth components studied by Oono & Suzuki (2020), which is closely related to normalized smoothness. Despite its importance, a theoretical understanding of normalized smoothness in GCNs with ReLU or leaky ReLU activation remains an open problem (Cai & Wang, 2020). Moreover, it is intriguing to ask whether such analysis could inspire new efficient algorithms to improve GCN performance.

1.1 Our Contribution

We aim to (1) understand how GCL smooths node features and (2) develop an efficient algorithm to let GCN and related models learn features with task-adaptive normalized smoothness (Cai & Wang, 2020) to improve node classification. Our main contributions are:

- We prove that there is a high-dimensional sphere underlying the input and output vectors of ReLU or leaky ReLU activation. This geometric characterization not only implies theories in (Oono & Suzuki, 2020; Cai & Wang, 2020), but also informs that adjusting the projection of input vectors onto the eigenspace \mathcal{M} can alter the smoothness of the output vectors (cf. Section 3).
- We show that both ReLU and leaky ReLU activation reduce the distance of features to the eigenspace \mathcal{M} , effectively smoothing node features without accounting for magnitude. However, when normalized by the feature magnitude, these activations can increase, decrease, or preserve smoothness—highlighting the nuanced behavior of normalized smoothness (cf. Sections 3 and 4).
- Inspired by our established geometric relationship between the input and output of ReLU or leaky ReLU activation, we study how modifying the projection of input features onto the eigenspace \mathcal{M} affects both normalized and unnormalized smoothness of the output vectors. We show that the output’s distance to \mathcal{M} never exceeds that of the input, regardless of how the projection is adjusted. In contrast, normalized smoothness can be effectively tuned to any desired level by controlling this projection (cf. Section 4).
- Based on our theoretical insights, we introduce a new smoothness control term (SCT) that enables GCNs and related models to learn node features with task-adaptive normalized smoothness, improving node classification. We validate SCT’s effectiveness across both homophilic and heterophilic graphs using several representative GCN-style models (cf. Sections 5 and 6).

To the best of our knowledge, this is the first comprehensive study of how ReLU and leaky ReLU activations affect the smoothness of node features—both in normalized and unnormalized forms.

1.2 Additional Related Works

Another line of related work focuses on controlling node feature smoothness to improve GCN performance. For example, Zhao & Akoglu (2020) introduce a normalization layer to prevent feature homogenization, while Zhou et al. (2021) constrain the Dirichlet energy without accounting for nonlinear activations. Although prior efforts have addressed over-smoothing and smoothness control, without considering activations, there remains a lack of theoretical analysis on how activation functions—particularly when accounting for feature magnitude—affect the smoothness of node features. A recent work (Di Giovanni et al., 2022) studies the normalized Dirichlet energy in GCNs with residual connections, demonstrating how the spectrum of the weight matrix affects smoothness. However, it does not offer an efficient method for controlling normalized smoothness.

1.3 Notation and Organization

Notation. We denote the ℓ_2 -norm of a vector \mathbf{u} as $\|\mathbf{u}\|$. For vectors \mathbf{u} and \mathbf{v} , we denote $\langle \mathbf{u}, \mathbf{v} \rangle$, $\mathbf{u} \odot \mathbf{v}$, and $\mathbf{u} \otimes \mathbf{v}$ as their inner, Hadamard, and Kronecker product, respectively. For a matrix \mathbf{A} , we denote its (i, j) th entry, transpose, and inverse as A_{ij} , \mathbf{A}^\top , and \mathbf{A}^{-1} , respectively. We denote the trace of matrix \mathbf{A} as $\text{Trace}(\mathbf{A}) = \sum_{i=1}^n A_{ii}$. For two matrices \mathbf{A} and \mathbf{B} , we denote their Frobenius inner product as $\langle \mathbf{A}, \mathbf{B} \rangle_F := \text{Trace}(\mathbf{A}\mathbf{B}^\top)$ and the Frobenius norm of matrix \mathbf{A} as $\|\mathbf{A}\|_F := \sqrt{\langle \mathbf{A}, \mathbf{A} \rangle}$.

Organization. We provide preliminaries and a review of some related results in Section 2. In Section 3, we establish a geometric characterization of how ReLU and leaky ReLU activations affect the smoothness of their input vectors. We study the smoothness of each dimension of node features and take their magnitude into account in Section 4. Our proposed SCT is presented in Section 5. We verify the efficacy of the proposed SCT for improving node classification in Section 6. Technical proofs and additional experimental details are provided in the appendix.

2 Preliminaries and Existing Results

By spectral graph theory (Chung, 1997), the eigenvalues of the normalized adjacency matrix \mathbf{G} in equation 1 can be ordered as $1 = \lambda_1 = \dots = \lambda_m > \lambda_{m+1} \geq \dots \geq \lambda_n > -1$, where m is the number of connected components in the graph G . Accordingly, the vertex set $V = \{v_k\}_{k=1}^n$ can be partitioned into m connected components V_1, \dots, V_m . Let $\mathbf{u}_i = (\mathbf{1}_{\{v_k \in V_i\}})_{1 \leq k \leq n}$ denote the indicator vector for component V_i , where the k -th entry is 1 if $v_k \in V_i$, and 0 otherwise. Let \mathbf{e}_i be the eigenvector corresponding to λ_i ; then $\{\mathbf{e}_i\}_{i=1}^n$ forms an orthonormal basis of \mathbb{R}^n . The eigenspace \mathcal{M} is spanned by $\{\mathbf{e}_i\}_{i=1}^m$, and its orthogonal complement \mathcal{M}^\perp by $\{\mathbf{e}_i\}_{i=m+1}^n$. Oono & Suzuki (2020) relate the indicator vectors \mathbf{u}_i to the eigenspace \mathcal{M} , specifically showing that:

Proposition 2.1 (Oono & Suzuki (2020)). *All eigenvalues of matrix \mathbf{G} lie in the interval $(-1, 1]$. Moreover, the nonnegative vectors $\{\tilde{\mathbf{D}}^{\frac{1}{2}}\mathbf{u}_i / \|\tilde{\mathbf{D}}^{\frac{1}{2}}\mathbf{u}_i\|\}_{1 \leq i \leq m}$ form an orthonormal basis of \mathcal{M} .*

For any matrix $\mathbf{H} := [\mathbf{h}_1, \dots, \mathbf{h}_n] \in \mathbb{R}^{d \times n}$, we have the decomposition

$$\mathbf{H} = \mathbf{H}_{\mathcal{M}} + \mathbf{H}_{\mathcal{M}^\perp},$$

where $\mathbf{H}_{\mathcal{M}} = \sum_{i=1}^m \mathbf{H}\mathbf{e}_i\mathbf{e}_i^\top$ and $\mathbf{H}_{\mathcal{M}^\perp} = \sum_{i=m+1}^n \mathbf{H}\mathbf{e}_i\mathbf{e}_i^\top$, satisfying

$$\langle \mathbf{H}_{\mathcal{M}}, \mathbf{H}_{\mathcal{M}^\perp} \rangle_F = \text{Trace} \left(\sum_{i=1}^m \mathbf{H}\mathbf{e}_i\mathbf{e}_i^\top \left(\sum_{j=m+1}^n \mathbf{H}\mathbf{e}_j\mathbf{e}_j^\top \right)^\top \right) = 0.$$

This implies that $\|\mathbf{H}\|_F^2 = \|\mathbf{H}_{\mathcal{M}}\|_F^2 + \|\mathbf{H}_{\mathcal{M}^\perp}\|_F^2$.

2.1 Existing Smoothness Notions of Node Features

Distance to the eigenspace \mathcal{M} . Oono & Suzuki (2020) study the smoothness of \mathbf{H} using their distance to the eigenspace \mathcal{M} as an unnormalized smoothness notion, i.e.,

Definition 2.2 (Oono & Suzuki (2020)). Let $\mathbb{R}^d \otimes \mathcal{M}$ be the subspace of $\mathbb{R}^{d \times n}$ consisting of the sum $\sum_{i=1}^m \mathbf{w}_i \otimes \mathbf{e}_i$, where $\mathbf{w}_i \in \mathbb{R}^d$ and $\{\mathbf{e}_i\}_{i=1}^m$ is an orthonormal basis of the eigenspace \mathcal{M} . Then we define $\|\mathbf{H}\|_{\mathcal{M}^\perp}$ —the distance of \mathbf{H} to \mathcal{M} —as follows:

$$\|\mathbf{H}\|_{\mathcal{M}^\perp} := \inf_{\mathbf{Y} \in \mathbb{R}^d \otimes \mathcal{M}} \|\mathbf{H} - \mathbf{Y}\|_F = \left\| \mathbf{H} - \sum_{i=1}^m \mathbf{H} \mathbf{e}_i \mathbf{e}_i^\top \right\|_F.$$

Since $\mathbf{H} = \mathbf{H}_{\mathcal{M}} + \mathbf{H}_{\mathcal{M}^\perp}$, then we have

$$\|\mathbf{H}\|_{\mathcal{M}^\perp} = \inf_{\mathbf{Y} \in \mathbb{R}^d \otimes \mathcal{M}} \|\mathbf{H} - \mathbf{Y}\|_F = \|\mathbf{H} - \mathbf{H}_{\mathcal{M}}\|_F = \|\mathbf{H}_{\mathcal{M}^\perp}\|_F. \quad (3)$$

Dirichlet energy. Cai & Wang (2020) study the unnormalized smoothness of node features using the following Dirichlet energy:

Definition 2.3 (Cai & Wang (2020)). Let $\tilde{\Delta} = \mathbf{I} - \mathbf{G}$ be the (augmented) normalized Laplacian, then the Dirichlet energy $\|\mathbf{H}\|_E$ of node features \mathbf{H} is defined by $\|\mathbf{H}\|_E^2 := \text{Trace}(\mathbf{H} \tilde{\Delta} \mathbf{H}^\top)$.

Normalized Dirichlet energy. Cai & Wang (2020) point out that the smoothness of \mathbf{H} should be measured using the normalized Dirichlet energy, defined as $\text{Trace}(\mathbf{H} \tilde{\Delta} \mathbf{H}^\top) / \|\mathbf{H}\|_F^2$. Normalizing the measurement can effectively mitigate biases result from these different scales.

2.2 Two Existing Theories of Over-smoothing

Let $\{\lambda_i\}_{i=1}^n$ be the eigenvalues of \mathbf{G} , and define $\lambda = \max\{|\lambda_i| \mid \lambda_i < 1\}$ as the second-largest magnitude. Let s_l denote the largest singular value of the weight matrix \mathbf{W}^l . Oono & Suzuki (2020) show that under the ReLU activation, GCL satisfies $\|\mathbf{H}^l\|_{\mathcal{M}^\perp} \leq s_l \lambda \|\mathbf{H}^{l-1}\|_{\mathcal{M}^\perp}$, implying $\|\mathbf{H}^l\|_{\mathcal{M}^\perp} \rightarrow 0$ as $l \rightarrow \infty$ if $s_l \lambda < 1$. This convergence to the eigenspace \mathcal{M} leads to over-smoothing. A key step in their analysis is the inequality $\|\sigma(\mathbf{Z})\|_{\mathcal{M}^\perp} \leq \|\mathbf{Z}\|_{\mathcal{M}^\perp}$ for any matrix \mathbf{Z} when σ is the ReLU activation, showing that the ReLU activation reduces the distance to \mathcal{M} . However, extending this result to other activations, including leaky ReLU, remains challenging (Oono & Suzuki, 2020).

Instead of analyzing $\|\mathbf{H}\|_{\mathcal{M}^\perp}$, Cai & Wang (2020) prove that under the GCL with ReLU or leaky ReLU activation, the normalized Dirichlet energy satisfies $\|\mathbf{H}^l\|_E \leq s_l \lambda \|\mathbf{H}^{l-1}\|_E$, implying $\|\mathbf{H}^l\|_E \rightarrow 0$ as $l \rightarrow \infty$ and thus over-smoothing. Their proof applies to both activation functions by establishing the inequality $\|\sigma(\mathbf{Z})\|_E \leq \|\mathbf{Z}\|_E$ for any matrix \mathbf{Z} .

3 Effects of Activation Functions: A Geometric Characterization

In this section, we present the geometric relationship between the input and output vectors of ReLU or leaky ReLU activation. For all subsequent analyses, we adopt $\|\mathbf{H}\|_{\mathcal{M}^\perp}$ as our unnormalized smoothness measure, noting its equivalence to $\|\mathbf{H}\|_E$ as a seminorm. Specifically, we have

Proposition 3.1. $\|\mathbf{H}\|_{\mathcal{M}^\perp}$ and $\|\mathbf{H}\|_E$ are two equivalent seminorms, i.e., there exist two constants $\alpha, \beta > 0$ s.t. $\alpha \|\mathbf{H}\|_{\mathcal{M}^\perp} \leq \|\mathbf{H}\|_E \leq \beta \|\mathbf{H}\|_{\mathcal{M}^\perp}$, for any $\mathbf{H} \in \mathbb{R}^{d \times n}$.

ReLU. Let $\sigma(x) = \max\{x, 0\}$ be the ReLU activation. Our first main result is that there is a high-dimensional sphere underlying the input and output vectors of the ReLU activation. More precisely, we have the following result:

Proposition 3.2 (ReLU). For any $\mathbf{Z} = \mathbf{Z}_{\mathcal{M}} + \mathbf{Z}_{\mathcal{M}^\perp} \in \mathbb{R}^{d \times n}$, let $\mathbf{H} = \sigma(\mathbf{Z}) = \mathbf{H}_{\mathcal{M}} + \mathbf{H}_{\mathcal{M}^\perp}$. Then $\mathbf{H}_{\mathcal{M}^\perp}$ lies on the high-dimensional sphere centered at $\mathbf{Z}_{\mathcal{M}^\perp}/2$ with radius

$$r := \left(\|\mathbf{Z}_{\mathcal{M}^\perp}/2\|_F^2 - \langle \mathbf{H}_{\mathcal{M}}, \mathbf{H}_{\mathcal{M}} - \mathbf{Z}_{\mathcal{M}} \rangle_F \right)^{1/2}.$$

In particular, $\mathbf{H}_{\mathcal{M}^\perp}$ lies inside the ball centered at $\mathbf{Z}_{\mathcal{M}^\perp}/2$ with radius $\|\mathbf{Z}_{\mathcal{M}^\perp}/2\|_F$ and hence

$$\|\mathbf{H}\|_{\mathcal{M}^\perp} \leq \|\mathbf{Z}\|_{\mathcal{M}^\perp}.$$

Leaky ReLU. For leaky ReLU activation $\sigma_a(x) = \max\{x, ax\}$ for $0 < a < 1$, we have

Proposition 3.3 (Leaky ReLU). *For any $\mathbf{Z} = \mathbf{Z}_{\mathcal{M}} + \mathbf{Z}_{\mathcal{M}^\perp} \in \mathbb{R}^{d \times n}$, let $\mathbf{H} = \sigma_a(\mathbf{Z}) = \mathbf{H}_{\mathcal{M}} + \mathbf{H}_{\mathcal{M}^\perp}$. Then $\mathbf{H}_{\mathcal{M}^\perp}$ lies on the sphere centered at $(1+a)\mathbf{Z}_{\mathcal{M}^\perp}/2$ with radius*

$$r_a := \left(\|(1-a)\mathbf{Z}_{\mathcal{M}^\perp}/2\|_F^2 - \langle \mathbf{H}_{\mathcal{M}} - \mathbf{Z}_{\mathcal{M}}, \mathbf{H}_{\mathcal{M}} - a\mathbf{Z}_{\mathcal{M}} \rangle_F \right)^{1/2}.$$

In particular, $\mathbf{H}_{\mathcal{M}^\perp}$ lies inside the ball centered at $(1+a)\mathbf{Z}_{\mathcal{M}^\perp}/2$ with radius $\|(1-a)\mathbf{Z}_{\mathcal{M}^\perp}/2\|_F$ and hence we see that

$$a\|\mathbf{Z}\|_{\mathcal{M}^\perp} \leq \|\mathbf{H}\|_{\mathcal{M}^\perp} \leq \|\mathbf{Z}\|_{\mathcal{M}^\perp}.$$

3.1 Implications of the Above Geometric Characterizations

Propositions 3.2 and 3.3 show that the unnormalized smoothness $\|\mathbf{H}_{\mathcal{M}^\perp}\|_F = \|\mathbf{H}\|_{\mathcal{M}^\perp}$ depends on the center and radii r or r_a . Given $\mathbf{Z}_{\mathcal{M}^\perp}$, the center of the spheres remains unchanged, and their radii r and r_a are only affected by changes in $\mathbf{Z}_{\mathcal{M}}$. This motivates our study of **how variations in $\mathbf{Z}_{\mathcal{M}}$ influence the unnormalized smoothness of node features**.

Furthermore, Propositions 3.2 and 3.3 imply both ReLU and leaky ReLU activation reduce the distance of node features to \mathcal{M} , i.e., $\|\mathbf{H}\|_{\mathcal{M}^\perp} \leq \|\mathbf{Z}\|_{\mathcal{M}^\perp}$. Moreover, this inequality is independent of $\mathbf{Z}_{\mathcal{M}}$; consider two node features $\mathbf{Z}, \mathbf{Z}' \in \mathbb{R}^{d \times n}$ s.t. $\mathbf{Z}_{\mathcal{M}^\perp} = \mathbf{Z}'_{\mathcal{M}^\perp}$ but $\mathbf{Z}_{\mathcal{M}} \neq \mathbf{Z}'_{\mathcal{M}}$. Let \mathbf{H} and \mathbf{H}' be the output of \mathbf{Z} and \mathbf{Z}' via ReLU or leaky ReLU activation, respectively. Then we have $\|\mathbf{H}\|_{\mathcal{M}^\perp} \leq \|\mathbf{Z}\|_{\mathcal{M}^\perp}$ and $\|\mathbf{H}'\|_{\mathcal{M}^\perp} \leq \|\mathbf{Z}'\|_{\mathcal{M}^\perp}$. Since $\mathbf{Z}_{\mathcal{M}^\perp} = \mathbf{Z}'_{\mathcal{M}^\perp}$, we deduce that $\|\mathbf{H}'\|_{\mathcal{M}^\perp} \leq \|\mathbf{Z}\|_{\mathcal{M}^\perp}$. In other words, when $\mathbf{Z}_{\mathcal{M}^\perp} = \mathbf{Z}'_{\mathcal{M}^\perp}$, changing $\mathbf{Z}_{\mathcal{M}}$ to $\mathbf{Z}'_{\mathcal{M}}$ can change the unnormalized smoothness of the output features but cannot change the fact that both ReLU and leaky ReLU activations smooth node features; we demonstrate this result in Fig. 1 a). Notice that without considering the nonlinear activation function, changing $\mathbf{Z}_{\mathcal{M}}$ does not affect the unnormalized smoothness of node features measured by $\|\mathbf{H}\|_{\mathcal{M}^\perp}$.

In contrast, **normalized smoothness can decrease when $\mathbf{Z}_{\mathcal{M}}$ is adjusted, potentially leading to a less smooth output vector**; see the next section for details.

4 Impact of Adjusting $\mathbf{Z}_{\mathcal{M}}$ on Output Smoothness

Throughout this section, let \mathbf{Z} and \mathbf{H} denote the input and output of ReLU or leaky ReLU activation. Traditional smoothness measures—based on distance to the eigenspace \mathcal{M} and Dirichlet energy—do not account for the relative magnitudes of feature dimensions. As noted by Cai & Wang (2020), analyzing the normalized smoothness $\|\mathbf{Z}\|_E/\|\mathbf{Z}\|_F$ remains an open challenging. These measures aggregate smoothness across all dimensions, so when certain dimensions dominate in magnitude, they disproportionately influence the overall smoothness.

Building on the discussion in Section 3.1, we examine how adjusting $\mathbf{Z}_{\mathcal{M}}$ differently affects normalized and unnormalized smoothness. For simplicity, we assume the graph is connected ($m = 1$); the results extend naturally to graphs with multiple components. Given the equivalence between the seminorms $\|\cdot\|_{\mathcal{M}}$ and $\|\cdot\|_E$, we introduce the following definition of dimension-wise normalized smoothness for node features:

Definition 4.1. Let $\mathbf{Z} \in \mathbb{R}^{d \times n}$ be the feature matrix over n graph nodes, where $\mathbf{z}^{(i)} \in \mathbb{R}^n$ denotes the i^{th} row of \mathbf{Z} , i.e., the i^{th} feature dimension across all nodes. We define the normalized smoothness of $\mathbf{z}^{(i)}$ as:

$$s(\mathbf{z}^{(i)}) := \|\mathbf{z}_{\mathcal{M}}^{(i)}\|/\|\mathbf{z}^{(i)}\|,$$

where we set $s(\mathbf{z}^{(i)}) = 1$ when $\mathbf{z}^{(i)} = \mathbf{0}$.

Remark 4.2. Notice that the normalized smoothness $s(\mathbf{z}^{(i)}) = \|\mathbf{z}_{\mathcal{M}}^{(i)}\|/\|\mathbf{z}^{(i)}\|$ is related to the ratio between the smooth and non-smooth components of node features $\|\mathbf{z}_{\mathcal{M}}^{(i)}\|/\|\mathbf{z}_{\mathcal{M}^\perp}^{(i)}\|$.

The graph is connected implies that

$$\mathbf{z}_{\mathcal{M}}^{(i)} = \langle \mathbf{z}^{(i)}, \mathbf{e}_1 \rangle \mathbf{e}_1 \implies \|\mathbf{z}_{\mathcal{M}}^{(i)}\| = |\langle \mathbf{z}^{(i)}, \mathbf{e}_1 \rangle|.$$

For clarity, we drop the index and write \mathbf{z} for $\mathbf{z}^{(i)}$ and \mathbf{e} for \mathbf{e}_1 , the unique eigenvector of \mathbf{G} associated with eigenvalue 1. Moreover, we have

$$s(\mathbf{z}) = \frac{\|\mathbf{z}_{\mathcal{M}}\|}{\|\mathbf{z}\|} = \frac{|\langle \mathbf{z}, \mathbf{e} \rangle|}{\|\mathbf{z}\|} = \frac{|\langle \mathbf{z}, \mathbf{e} \rangle|}{\|\mathbf{z}\| \cdot \|\mathbf{e}\|} \implies 0 \leq s(\mathbf{z}) \leq 1. \quad (4)$$

It is evident that the larger $s(\mathbf{z})$ is, the smoother the node feature \mathbf{z} is. In fact, we have

$$s(\mathbf{z})^2 + \left(\frac{\|\mathbf{z}\|_{\mathcal{M}^\perp}}{\|\mathbf{z}\|} \right)^2 = \frac{\|\mathbf{z}_{\mathcal{M}}\|^2}{\|\mathbf{z}\|^2} + \frac{\|\mathbf{z}_{\mathcal{M}^\perp}\|^2}{\|\mathbf{z}\|^2} = 1,$$

where $\|\mathbf{z}\|_{\mathcal{M}^\perp}/\|\mathbf{z}\|$ decreases as $s(\mathbf{z})$ increases.

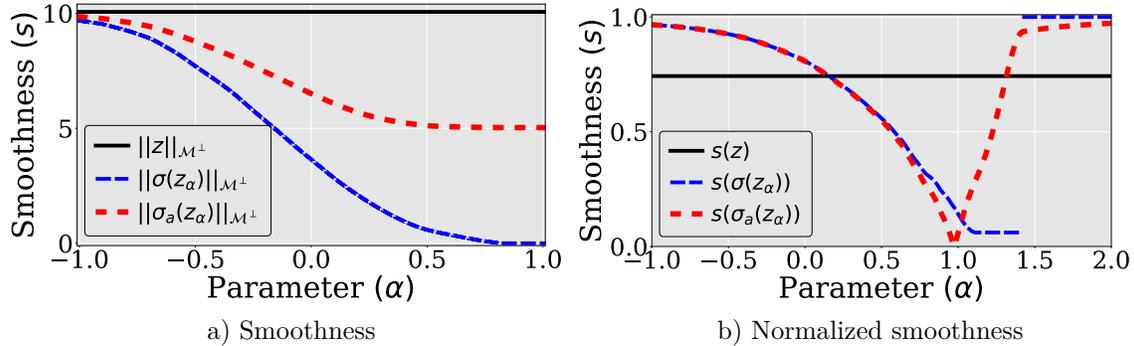


Figure 1: Contrasting the effects of varying the parameter α on the smoothness and normalized smoothness of the output features $\sigma(\mathbf{z}_\alpha)$ and $\sigma_a(\mathbf{z}_\alpha)$. The discontinuity in $s(\sigma(\mathbf{z}_\alpha))$ shown in panel b) arises from the definition of normalized smoothness. Note that $s(\mathbf{z}) = 1$ when $\mathbf{z} = \mathbf{0}$, and $\sigma(\mathbf{z}_\alpha)$ can become $\mathbf{0}$ for sufficiently large α .

To discuss how the smoothness $s(\mathbf{h}) = s(\sigma(\mathbf{z}))$ or $s(\sigma_a(\mathbf{z}))$ can be adjusted by changing $\mathbf{z}_{\mathcal{M}}$, we define

$$\mathbf{z}(\alpha) := \mathbf{z} - \alpha \mathbf{e}.$$

It is clear that $\mathbf{z}(\alpha)_{\mathcal{M}^\perp} = \mathbf{z}_{\mathcal{M}^\perp}$ and $\mathbf{z}(\alpha)_{\mathcal{M}} = \mathbf{z}_{\mathcal{M}} - \alpha \mathbf{e}$. We see that α only alters $\mathbf{z}_{\mathcal{M}}$ while preserves $\mathbf{z}_{\mathcal{M}^\perp}$. Moreover, we have

$$s(\mathbf{z}(\alpha)) = \sqrt{1 - \frac{\|\mathbf{z}(\alpha)_{\mathcal{M}^\perp}\|^2}{\|\mathbf{z}(\alpha)\|^2}} = \sqrt{1 - \frac{\|\mathbf{z}_{\mathcal{M}^\perp}\|^2}{\|\mathbf{z}(\alpha)\|^2}}.$$

It follows that $s(\mathbf{z}(\alpha)) = 1$ if and only if $\mathbf{z}_{\mathcal{M}^\perp} = \mathbf{0}$, showing that when $\mathbf{z}_{\mathcal{M}^\perp} = \mathbf{0}$, the vector \mathbf{z} is the smoothest one.

4.1 The Disparate Effects of α on $\|\cdot\|_{\mathcal{M}^\perp}$ and $s(\cdot)$: Numerical Studies

To explore how the smoothness measures vary with α , we perform some empirical studies in this subsection.

Let $\mathbf{z}_\alpha := \mathbf{z}(\alpha) = \mathbf{z} - \alpha \mathbf{e}$, where $\mathbf{z} \in \mathbb{R}^{100}$ is a randomly initialized node feature vector (uniformly sampled from $[-1.5, 1.5]$), and \mathbf{e} is computed as $\mathbf{e} = \tilde{\mathbf{D}}^{1/2} \mathbf{u} / \|\tilde{\mathbf{D}}^{1/2} \mathbf{u}\|$ per Proposition 2.1, with \mathbf{u} being the all-ones vector and $\tilde{\mathbf{D}}$ the augmented degree matrix. We use a connected synthetic graph with 100 nodes, each assigned a random degree between 2 and 10.

We evaluate both unnormalized smoothness $\|\sigma(\mathbf{z}_\alpha)\|_{\mathcal{M}^\perp}$, $\|\sigma_a(\mathbf{z}_\alpha)\|_{\mathcal{M}^\perp}$ and normalized smoothness $s(\sigma(\mathbf{z}_\alpha))$, $s(\sigma_a(\mathbf{z}_\alpha))$ across $\alpha \in [-1.5, 1.5]$. As shown in Fig. 1 a), the output smoothness is consistently no greater

than the input smoothness $\|\mathbf{z}\|_{\mathcal{M}^\perp}$, confirming the result from Section 3.1: modifying $\mathbf{z}_{\mathcal{M}}$ does not affect this inequality. However, an important observation is that **adjusting the eigenspace projection can still alter the output smoothness**, even though the input smoothness remains unchanged.

In contrast, Fig. 1 b) shows that the normalized smoothness $s(\sigma(\mathbf{z}(\alpha)))$ and $s(\sigma_a(\mathbf{z}(\alpha)))$ can be reduced below $s(\mathbf{z})$ by adjusting α . This highlights how modifying the projection onto \mathcal{M} can influence normalized smoothness, even when unnormalized smoothness remains bounded.

4.2 The Smooth Effects of ReLU and Leaky ReLU: Theoretical Analysis

In this subsection, we provide theoretical insights into the empirical findings shown in Fig. 1, focusing on how the smoothness of $\sigma(\mathbf{z}(\alpha))$ and $\sigma_a(\mathbf{z}(\alpha))$ varies with α . If $\mathbf{z}_{\mathcal{M}^\perp} = \mathbf{0}$, Propositions 3.2 and 3.3 imply that $\|\sigma(\mathbf{z}(\alpha))\|_{\mathcal{M}^\perp}$ and $\|\sigma_a(\mathbf{z}(\alpha))\|_{\mathcal{M}^\perp}$ vanish, leading to $s(\sigma(\mathbf{z}(\alpha))) = 1$ for all α . Therefore, we assume $\mathbf{z}_{\mathcal{M}^\perp} \neq \mathbf{0}$ in the following analysis.

Proposition 4.3 (ReLU). *Suppose $\mathbf{z}_{\mathcal{M}^\perp} \neq \mathbf{0}$. Let $\mathbf{h}(\alpha) = \sigma(\mathbf{z}(\alpha))$ with σ being the ReLU activation, then we have*

$$\min_{\alpha} s(\mathbf{h}(\alpha)) = \sqrt{\frac{\sum_{x_i=\max \mathbf{x}} d_i}{\sum_{j=1}^n d_j}}, \text{ and } \max_{\alpha} s(\mathbf{h}(\alpha)) = 1,$$

where $\mathbf{x} := \tilde{\mathbf{D}}^{-\frac{1}{2}} \mathbf{z}$, $\max \mathbf{x} = \max_{1 \leq i \leq n} x_i$, and $\tilde{\mathbf{D}}$ is the augmented degree matrix, whose diagonal entries are d_1, d_2, \dots, d_n . In particular, the normalized smoothness $s(\mathbf{h}(\alpha))$ is monotone increasing as α decreases whenever $\alpha < \|\tilde{\mathbf{D}}^{\frac{1}{2}} \mathbf{u}_n\| \max \mathbf{x}$ and it has range $[\min_{\alpha} s(\mathbf{h}(\alpha)), 1]$.

Proposition 4.4 (Leaky ReLU). *Assume $\mathbf{z}_{\mathcal{M}^\perp} \neq \mathbf{0}$ and let $\mathbf{h}(\alpha) = \sigma_a(\mathbf{z}(\alpha))$, where σ_a is the leaky ReLU activation. Then (1) $\min_{\alpha} s(\mathbf{h}(\alpha)) = 0$, and (2) $\sup_{\alpha} s(\mathbf{h}(\alpha)) = 1$ and $s(\mathbf{h}(\alpha))$ has range $[0, 1]$.*

Proposition 4.4 also holds for other variants of ReLU, e.g., ELU¹ and SELU²; see Appendix B. We summarize Propositions 3.2, 3.3, 4.3, and 4.4 in the following corollary, which qualitatively explains the empirical results in Fig. 1.

Corollary 4.5. *Suppose $\mathbf{z}_{\mathcal{M}^\perp} \neq \mathbf{0}$. Let $\mathbf{h}(\alpha) = \sigma(\mathbf{z}(\alpha))$ or $\sigma_a(\mathbf{z}(\alpha))$ with σ being the ReLU activation and σ_a being the leaky ReLU activation. Then $\|\mathbf{z}\|_{\mathcal{M}^\perp} \geq \|\mathbf{h}(\alpha)\|_{\mathcal{M}^\perp}, \forall \alpha \in \mathbb{R}$; however, $s(\mathbf{h}(\alpha))$ can be smaller than, larger than, or equal to $s(\mathbf{z})$ for different α .*

Propositions 4.3, 4.4, and Corollary 4.5, provide theoretical support for the empirical results in Fig. 1. These results show that modifying the projection $\mathbf{z}_{\mathcal{M}}$ can influence both the unnormalized and normalized smoothness of the output $\mathbf{h} = \sigma(\mathbf{z})$ or $\sigma_a(\mathbf{z})$. In particular, the normalized smoothness can be tuned to any value within the ranges established in the propositions. This insight opens the door to designing algorithms that control feature smoothness to enhance GCN performance—a direction we explore in the next section.

5 Controlling Smoothness of Node Features

While the ideal smoothness level for node features in a given classification task is unknown, our theory shows that both normalized and unnormalized smoothness can be modulated by adjusting the input’s projection onto \mathcal{M} . Motivated by this, we propose the following learnable smoothness control term to dynamically regulate the smoothness of node features:

$$\mathbf{B}_{\alpha}^l = \sum_{i=1}^m \alpha_i^l \mathbf{e}_i^{\top}, \quad (5)$$

where l is the layer index, $\{\mathbf{e}_i\}_{i=1}^m$ is the orthonormal basis of the eigenspace \mathcal{M} , and $\boldsymbol{\alpha}^l := \{\alpha_i^l\}_{i=1}^m$ is a collection of learnable vectors with $\alpha_i^l \in \mathbb{R}^d$ being approximated by a multi-layer perceptron (MLP). The

¹ELU: $f(x) = \max(x, 0) + \min(0, a \cdot (e^x - 1))$ where $a > 0$.

²SELU: $f(x) = c(\max(x, 0) + \min(0, a \cdot (e^x - 1)))$ where $a, c > 0$.

detailed configuration of α_i^l will be specified in each experiment later. One can see that \mathbf{B}_α^l always lies in $\mathbb{R}^d \otimes \mathcal{M}$. We integrate SCT into GCL, resulting in the following update equation:

$$\mathbf{H}^l = \sigma(\mathbf{W}^l \mathbf{H}^{l-1} \mathbf{G} + \mathbf{B}_\alpha^l). \quad (6)$$

We call the corresponding model GCN-SCT.

We now detail the proposed SCT. For each layer l , we introduce a learnable matrix $\mathbf{A}^l \in \mathbb{R}^{d \times m}$, where each column α_i^l corresponds to a basis vector in the eigenspace \mathcal{M} . Here, d is the feature dimension and m is the dimension of \mathcal{M} . We observe that SCT performs best when guided by degree pooling over graph subcomponents. Let $\mathbf{Q} := [\mathbf{e}_1, \dots, \mathbf{e}_m] \in \mathbb{R}^{n \times m}$ be the matrix of orthogonal basis vectors. Pooling is performed via $\mathbf{H}^l \mathbf{Q}$. For the first architecture, we let

$$\mathbf{A}^l = \mathbf{W} \odot (\mathbf{H}^l \mathbf{Q}),$$

where $\mathbf{W} \in \mathbb{R}^{d \times m}$ is learnable and performs pooling over \mathbf{H}^l using the eigenvectors \mathbf{Q} . The second architecture uses a residual connection with hyperparameter $\beta_l = \log(\theta/l + 1)$ and learnable matrices $\mathbf{W}_0, \mathbf{W}_1 \in \mathbb{R}^{d \times d}$ and the softmax function ϕ . Resulting in

$$\mathbf{A}^l = \phi(\mathbf{H}^l \mathbf{Q}) \odot (\beta_l \mathbf{W}_0 \mathbf{H}^0 \mathbf{Q} + (1 - \beta_l) \mathbf{W}_1 \mathbf{H}^l \mathbf{Q}).$$

In Section 6, we use the first architecture for GCN-SCT as GCN uses only \mathbf{H}^l information at each layer. We use the second architecture for GCNII-SCT and EGNN-SCT, which use both \mathbf{H}^0 and \mathbf{H}^l information at each layer. This design offers two key advantages: (1) it enables effective control of normalized smoothness, and (2) it is computationally efficient, relying only on eigenvectors associated with eigenvalue 1 of matrix \mathbf{G} .

Remark 5.1. Computing the basis of eigenspace \mathcal{M} (eigenvectors corresponding to eigenvalue 1 of \mathbf{G}) does not introduce substantial computational overhead. In particular, the basis of \mathcal{M} is given by the indicator functions of each connected component of the graph, and we can identify connected components for undirected graphs using disjoint set union (DSU) Galil & Italiano (1991) with linear time complexity with respect to the number of nodes.

5.1 Integrating SCT into Other GCN-style Models

In this subsection, we present two other representative applications of SCT. First, we apply SCT to GCNII (Chen et al., 2020b), a state-of-the-art (SOTA) GCN-style model (Chen et al., 2020b; Luan et al., 2022), to demonstrate that SCT can enhance node classification by improving feature smoothness. Second, we apply SCT to EGNN (Zhou et al., 2021), which originally controls feature smoothness via Dirichlet energy bounds under linear activation. We show that our new theoretical insights into activation functions, combined with SCT, can improve EGNN performance when nonlinear activations are considered.

GCNII. Each GCNII layer uses a skip connection to the initial layer \mathbf{H}^0 and is given as follows:

$$\mathbf{H}^l = \sigma(((1 - \alpha_l) \mathbf{H}^{l-1} \mathbf{G} + \alpha_l \mathbf{H}^0) ((1 - \beta_l) \mathbf{I} + \beta_l \mathbf{W}^l)),$$

where $\alpha_l, \beta_l \in (0, 1)$ are learnable. We integrate SCT \mathbf{B}_α^l into GCNII, resulting in the following GCNII-SCT layers:

$$\mathbf{H}^l = \sigma(((1 - \alpha_l) \mathbf{H}^{l-1} \mathbf{G} + \alpha_l \mathbf{H}^0) ((1 - \beta_l) \mathbf{I} + \beta_l \mathbf{W}^l) + \mathbf{B}_\alpha^l).$$

We call the resulting model GCNII-SCT.

EGNN. EGNN (Zhou et al., 2021) controls the smoothness of node features by constraining the lower and upper bounds of the Dirichlet energy of node features without considering the nonlinear activation function. Each EGNN layer can be written as follows:

$$\mathbf{H}^l = \sigma(\mathbf{W}^l (c_1 \mathbf{H}^0 + c_2 \mathbf{H}^{l-1} + (1 - c_{\min}) \mathbf{H}^{l-1} \mathbf{G})), \quad (7)$$

where c_1, c_2 are learnable weights that satisfy $c_1 + c_2 = c_{\min}$ with c_{\min} being a hyperparameter. To constrain the Dirichlet energy, EGNN initializes trainable weights \mathbf{W}^l as diagonal matrices with prescribed singular

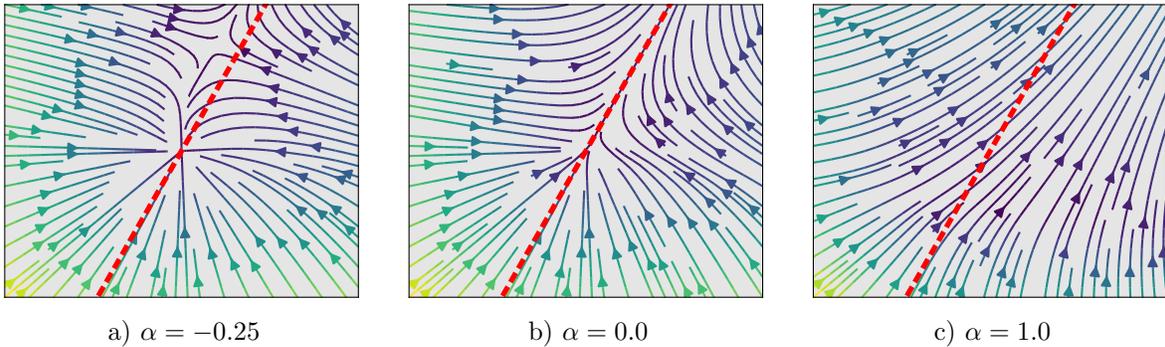


Figure 2: Node feature trajectories, with colored magnitude, for varying smoothness control parameter α . For classical GCN b), the node features converge to the eigenspace \mathcal{M} (red dashed line).

values and regularizes them to maintain orthogonality during training. Ignoring the activation function σ , the node features \mathbf{H}^l at layer l satisfy:

$$c_{\min}\|\mathbf{H}^0\|_E \leq \|\mathbf{H}^l\|_E \leq c_{\max}\|\mathbf{H}^0\|_E,$$

where c_{\max} is the square of the maximal singular value of the initialization of \mathbf{W}^1 . Similarly, we modify EGNN to result in the following EGNN-SCT layer:

$$\mathbf{H}^l = \sigma(\mathbf{W}^l((1 - c_{\min})\mathbf{H}^{l-1}\mathbf{G} + c_1\mathbf{H}^0 + c_2\mathbf{H}^{l-1}) + \mathbf{B}_\alpha^l),$$

where everything remains the same as the EGNN layer, except that we include our proposed SCT \mathbf{B}_α^l .

6 Experiments

We evaluate SCT on a range of benchmark node classification tasks to demonstrate its effectiveness in GCN-style models. Datasets include citation networks (Cora, Citeseer, PubMed, Coauthor-Physics, Ogbn-arxiv), web knowledge bases (Cornell, Texas, Wisconsin), and Wikipedia networks (Chameleon, Squirrel); see Appendix C.1 for details. We implement GCN (Kipf & Welling, 2017) and GCNII (Chen et al., 2020b) (without weight sharing) using PyTorch Geometric (PyG) (Fey & Lenssen, 2019), and use the official codebase for EGNN (Zhou et al., 2021)³.

6.1 Node Feature Trajectory

Following Oono & Suzuki (2020), we visualize the trajectory of node features for a simple graph with two connected nodes and 1D features. In this case, Eq. equation 6 simplifies to $\mathbf{h}^1 = \sigma(w\mathbf{h}^0\mathbf{G} + \mathbf{b}_\alpha)$, where $w = 1.2$, $\mathbf{h}^0, \mathbf{h}^1, \mathbf{b}_\alpha \in \mathbb{R}^2$, and $\mathbf{G} \in \mathbb{R}^{2 \times 2}$ with $\mathbf{G} = [0.592, 0.194; 0.194, 0.908]$, a positive definite matrix with largest eigenvalue 1. We sample 20 initial vectors \mathbf{h}^0 uniformly from $[-1, 1] \times [-1, 1]$. Figure 2 shows the resulting trajectories relative to the eigenspace \mathcal{M} (red dashed line). In panel a), some trajectories do not converge directly to \mathcal{M} ; in b), setting $\alpha = 0.0$ recovers GCL, and all trajectories converge to \mathcal{M} ; in c), large α (e.g., 1.0) causes significant initial deviation. These results show that α effectively controls feature trajectories.

6.2 Baseline Comparisons for Node Classification

6.2.1 Citation Networks

We compare the three GCN-style models from Section 5, with and without our proposed SCT, across varying depths in Table 1. For Cora, Citeseer, and PubMed, we use fixed splits from (Yang et al., 2016);

³<https://github.com/Kaixiong-Zhou/EGNN>

Layers	2	4	16	32
Cora				
GCN/GCN-SCT	81.1/ 82.9	80.4/ 82.8	64.9/ 71.4	60.3/ 67.2
GCNII/GCNII-SCT	82.2/ 83.8	82.6/ 84.3	84.6/ 84.8	85.4/ 85.5
EGNN/EGNN-SCT	83.2/ 84.1	84.2/ 84.5	85.4 /83.3	85.3 /82.0
Citeseer				
GCN/GCN-SCT	70.3 /69.9	67.6/ 67.7	18.3/ 55.4	25.0/ 51.0
GCNII/GCNII-SCT	68.2/ 72.8	68.9/ 72.8	72.9/ 73.8	73.4 / 73.4
EGNN/EGNN-SCT	72.0/ 73.1	71.9/ 72.0	72.4/ 72.6	72.3/ 72.9
PubMed				
GCN/GCN-SCT	79.0/ 79.8	76.5/ 78.4	40.9/ 76.1	22.4/ 77.0
GCNII/GCNII-SCT	78.2/ 79.7	78.8/ 80.1	80.2/ 80.7	79.8/ 80.7
EGNN/EGNN-SCT	79.2/ 79.8	79.5/ 80.4	80.1/ 80.3	80.0/ 80.4
Coauthor-Physics				
GCN/GCN-SCT	92.4/ 92.6 \pm 1.6	92.1/ 92.5 \pm 5.9	13.5/ 50.9 \pm 15.0	13.1/ 43.6 \pm 16.0
GCNII/GCNII-SCT	92.5/ 94.4 \pm 0.4	92.9/ 94.2 \pm 0.3	92.9/ 93.7 \pm 0.7	92.9/ 94.1 \pm 0.3
EGNN/EGNN-SCT	92.6/ 93.9 \pm 0.7	92.9/ 94.1 \pm 0.4	93.1/ 94.0 \pm 0.7	93.3/ 93.8 \pm 1.3
Ogbn-arxiv				
GCN/GCN-SCT	70.4/ 72.1 \pm 0.3	71.7/ 72.7 \pm 0.3	70.6/ 72.3 \pm 0.2	68.5/ 72.3 \pm 0.3
GCNII/GCNII-SCT	70.1/ 72.0 \pm 0.3	71.4/ 72.2 \pm 0.2	71.5/ 72.4 \pm 0.3	70.5/ 72.1 \pm 0.3
EGNN/EGNN-SCT	68.4/ 68.5 \pm 0.6	71.1/ 71.3 \pm 0.5	72.7/ 72.8 \pm 0.5	72.7 /72.3 \pm 0.5

Table 1: Model accuracies across varying depths on citation datasets. While GCN-SCT with 16 or 32 layers shows accuracy drops, this is due to vanishing gradients—not over-smoothing. For Cora, Citeseer, and PubMed, we follow Chen et al. (2020b) using fixed splits and a single forward pass, reporting only test accuracy. For Coauthor-Physics and Ogbn-arxiv, we use the splits from Zhou et al. (2021) and report both test accuracy and standard deviation. Baseline results are taken from Chen et al. (2020b); Zhou et al. (2021); standard deviations were not reported. (Unit: %)

for Coauthor-Physics and Ogbn-arxiv, we follow (Zhou et al., 2021). Dataset details are in Appendix C. Following (Chen et al., 2020b), we train using Adam (Kingma & Ba, 2014) with a single pass, 1500 max epochs, and 100-epoch patience. Hyperparameter grids are listed in Table 6, and we accelerate tuning via Bayesian meta-learning (Biewald, 2020), running 200 iterations per model.

Table 1 reports the best test accuracy (ReLU vs. leaky ReLU) for GCN, GCNII, and their SCT variants⁴. EGNN uses the SReLU activation as in (Zhou et al., 2021). Results show SCT consistently improves GCN and GCNII. EGNN-SCT (with ReLU or leaky ReLU activation) may underperform EGNN (with SReLU activation), due to activation choice. However, Appendix C.3 shows that EGNN-SCT with SReLU activation outperforms EGNN across all tasks. Since SReLU is a shifted ReLU, our theory still applies. Model sizes and runtimes are reported in Table 5 in the appendix.

Table 1 shows that GCN accuracy drops at depths 16 and 32, even with SCT. To investigate, we analyze the normalized smoothness of node features in GCN and GCN-SCT. Figure 3 visualizes heatmaps for Citeseer with 32 layers, where rows represent layers and columns represent feature dimensions. In GCN (Fig. 3 a), features from layers 14–32 exhibit smoothness near 1, indicating homogeneity. In contrast, GCN-SCT (Fig. 3 b) produces more diverse features. This suggests the performance drop in deep GCN-SCT is not due to over-smoothing. Unlike GCNII and EGNN, GCN-SCT lacks skip connections, which help mitigate vanishing gradients (He et al., 2016a;b). As shown in Appendix C.3, GCN and GCN-SCT suffer from vanishing gradients, while the other models do not.

⁴See Appendix C for full comparisons.

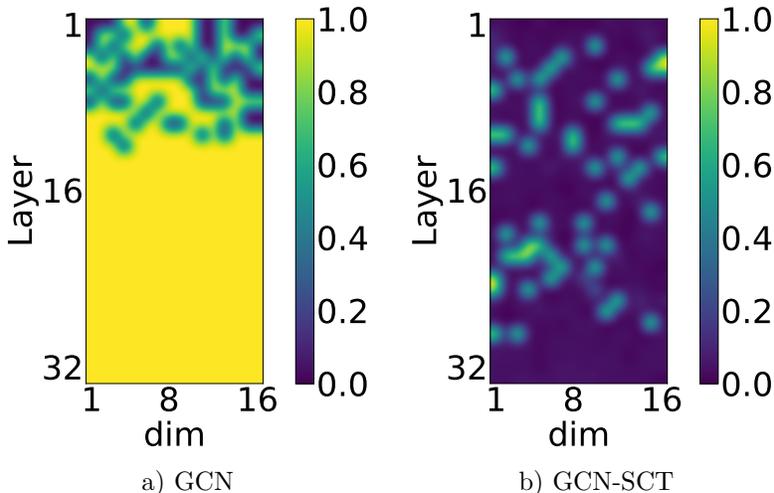


Figure 3: Normalized smoothness heatmaps for each feature dimension across layers in a) GCN and b) GCN-SCT on Citeseer (32 layers, 16 hidden dimensions). GCN features become fully smooth from layer 14 onward, while GCN-SCT maintains control over smoothness at all depths. Horizontal and vertical axes denote feature dimension and layer index, respectively.

6.2.2 Other Datasets

We further evaluate model performance using 10-fold cross-validation and fixed 48/32/20% splits, following (Pei et al., 2020). Tables 2 and 3 report accuracy and per-epoch runtime for GCN and GCNII (with and without SCT, using the leaky ReLU activation) on five heterophilic datasets: Cornell, Texas, Wisconsin, Chameleon, and Squirrel. EGNN is excluded, as these datasets were not considered in (Zhou et al., 2021). Baseline GCN and GCNII results are from (Chen et al., 2020b); all other models are tuned via Bayesian meta-learning to maximize validation accuracy. We report the best test accuracy over depths $\{2, 4, 8, 16, 32\}$.

SCT significantly improves accuracy and, in some cases, reduces runtime by achieving optimal performance at moderate depths. Table 9 (Section C.4) reports mean and standard deviation of test accuracy, while Table 10 shows that SCT adds minimal overhead at depth 8.

Cornell	Texas	Wisconsin	Chameleon	Squirrel
52.70/ 55.95	52.16/ 62.16	45.88/ 54.71	28.18/ 38.44	23.96/ 35.31
74.86/ 75.41	69.46/ 83.34	74.12/ 86.08	60.61/ 64.52	38.47/ 47.51

Table 2: Mean test accuracy on WebKB and WikipediaNetwork with fixed 48/32/20% splits. Rows: (1) GCN/GCN-SCT, (2) GCNII/GCNII-SCT. (Unit: %)

Cornell	Texas	Wisconsin	Chameleon	Squirrel
0.7/1.8	0.7/0.8	0.7/0.8	0.6/0.7	1.6/4.0
2.0/2.0	3.1/2.0	2.0/1.5	1.5/1.3	5.5/3.7

Table 3: Average epoch time on WebKB and WikipediaNetwork with fixed 48/32/20% splits. Rows: (1) GCN/GCN-SCT, (2) GCNII/GCNII-SCT. (Unit: $\times 10^{-2}$ second)

7 Concluding Remarks

We established a geometric understanding of how ReLU and leaky ReLU activation functions influence GCN feature smoothness. Our analysis of dimension-wise normalized smoothness shows that activation functions

not only smooth features but can also reduce or preserve their normalized smoothness. These insights guided the design of SCT, a new simple yet effective mechanism for controlling both normalized and unnormalized smoothness in GCN features. SCT offers theoretical guarantees for smoothness control in GCN-style models.

References

- Justin Baker, Qingsong Wang, Cory D Hauck, and Bao Wang. Implicit graph neural networks: A monotone operator viewpoint. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett (eds.), *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pp. 1521–1548. PMLR, 23–29 Jul 2023. URL <https://proceedings.mlr.press/v202/baker23a.html>.
- Justin M Baker, Qingsong Wang, Martin Berzins, Thomas Strohmer, and Bao Wang. Monotone operator theory-inspired message passing for learning long-range interaction on graphs. In *International Conference on Artificial Intelligence and Statistics*, pp. 2233–2241. PMLR, 2024.
- 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.
- Lukas Biewald. Experiment tracking with weights and biases, 2020. URL <https://www.wandb.com/>. Software available from wandb.com.
- Chen Cai and Yusu Wang. A note on over-smoothing for graph neural networks. *arXiv preprint arXiv:2006.13318*, 2020.
- Ben Chamberlain, James Rowbottom, Maria I Gorinova, Michael Bronstein, Stefan Webb, and Emanuele Rossi. Grand: Graph neural diffusion. In Marina Meila and Tong Zhang (eds.), *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pp. 1407–1418. PMLR, 18–24 Jul 2021. URL <https://proceedings.mlr.press/v139/chamberlain21a.html>.
- Deli Chen, Yankai Lin, Wei Li, Peng Li, Jie Zhou, and Xu Sun. Measuring and relieving the over-smoothing problem for graph neural networks from the topological view. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pp. 3438–3445, 2020a.
- Ming Chen, Zhewei Wei, Zengfeng Huang, Bolin Ding, and Yaliang Li. Simple and deep graph convolutional networks. In Hal Daumé III and Aarti Singh (eds.), *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pp. 1725–1735. PMLR, 13–18 Jul 2020b.
- Zhengdao Chen, Lisha Li, and Joan Bruna. Supervised community detection with line graph neural networks. In *International Conference on Learning Representations*, 2019.
- Fan RK Chung. *Spectral graph theory*, volume 92. American Mathematical Soc., 1997.
- Michaël Defferrard, Xavier Bresson, and Pierre Vandergheynst. Convolutional neural networks on graphs with fast localized spectral filtering. *Advances in neural information processing systems*, 29, 2016.
- Francesco Di Giovanni, James Rowbottom, Benjamin P Chamberlain, Thomas Markovich, and Michael M Bronstein. Understanding convolution on graphs via energies. *arXiv preprint arXiv:2206.10991*, 2022.
- Matthias Fey and Jan E. Lenssen. Fast graph representation learning with PyTorch Geometric. In *ICLR Workshop on Representation Learning on Graphs and Manifolds*, 2019.
- Zvi Galil and Giuseppe F Italiano. Data structures and algorithms for disjoint set union problems. *ACM Computing Surveys (CSUR)*, 23(3):319–344, 1991.
- Johannes Gasteiger, Aleksandar Bojchevski, and Stephan Günnemann. Combining neural networks with personalized pagerank for classification on graphs. In *International Conference on Learning Representations*, 2019.

- Justin Gilmer, Samuel S. Schoenholz, Patrick F. Riley, Oriol Vinyals, and George E. Dahl. Neural message passing for quantum chemistry. In *Proceedings of the 34th International Conference on Machine Learning - Volume 70*, ICML'17, pp. 1263–1272. JMLR.org, 2017.
- Fangda Gu, Heng Chang, Wenwu Zhu, Somayeh Sojoudi, and Laurent El Ghaoui. Implicit graph neural networks. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (eds.), *Advances in Neural Information Processing Systems*, volume 33, pp. 11984–11995. Curran Associates, Inc., 2020. URL <https://proceedings.neurips.cc/paper/2020/file/8b5c8441a8ff8e151b191c53c1842a38-Paper.pdf>.
- William Hamilton, Zhitao Ying, and Jure Leskovec. Inductive representation learning on large graphs. *Advances in neural information processing systems*, 30, 2017a.
- William L Hamilton. *Graph representation learning*. Morgan & Claypool Publishers, 2020.
- William L Hamilton, Rex Ying, and Jure Leskovec. Representation learning on graphs: Methods and applications. *arXiv preprint arXiv:1709.05584*, 2017b.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016a.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Identity mappings in deep residual networks. In *Computer Vision—ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part IV 14*, pp. 630–645. Springer, 2016b.
- Tatsuro Kawamoto, Masashi Tsubaki, and Tomoyuki Obuchi. Mean-field theory of graph neural networks in graph partitioning. *Advances in Neural Information Processing Systems*, 31, 2018.
- Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Thomas N. Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. In *International Conference on Learning Representations*, 2017.
- Qimai Li, Zhichao Han, and Xiao-Ming Wu. Deeper insights into graph convolutional networks for semi-supervised learning. In *Thirty-Second AAAI conference on artificial intelligence*, 2018.
- Sitao Luan, Chenqing Hua, Qincheng Lu, Jiaqi Zhu, Mingde Zhao, Shuyuan Zhang, Xiao-Wen Chang, and Doina Precup. Revisiting heterophily for graph neural networks. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho (eds.), *Advances in Neural Information Processing Systems*, 2022.
- Hoang Nt and Takanori Maehara. Revisiting graph neural networks: All we have is low-pass filters. *arXiv preprint arXiv:1905.09550*, 2019.
- Kenta Oono and Taiji Suzuki. Graph neural networks exponentially lose expressive power for node classification. In *International Conference on Learning Representations*, 2020.
- Hongbin Pei, Bingzhe Wei, Kevin Chen-Chuan Chang, Yu Lei, and Bo Yang. Geom-gcn: Geometric graph convolutional networks. In *International Conference on Learning Representations*, 2020.
- Yu Rong, Wenbing Huang, Tingyang Xu, and Junzhou Huang. Droppedge: Towards deep graph convolutional networks on node classification. In *International Conference on Learning Representations*, 2020.
- Michael Scholkemper, Xinyi Wu, Ali Jadbabaie, and Michael T Schaub. Residual connections and normalization can provably prevent oversmoothing in gnns. *arXiv preprint arXiv:2406.02997*, 2024.
- Matthew Thorpe, Tan Minh Nguyen, Hedi Xia, Thomas Strohmmer, Andrea Bertozzi, Stanley Osher, and Bao Wang. GRAND++: Graph neural diffusion with a source term. In *International Conference on Learning Representations*, 2022. URL <https://openreview.net/forum?id=EMxu-dzvJk>.
- Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua Bengio. Graph attention networks. In *International Conference on Learning Representations*, 2018.

- Felix Wu, Amauri Souza, Tianyi Zhang, Christopher Fifty, Tao Yu, and Kilian Weinberger. Simplifying graph convolutional networks. In Kamalika Chaudhuri and Ruslan Salakhutdinov (eds.), *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pp. 6861–6871. PMLR, 09–15 Jun 2019.
- Xinyi Wu, Amir Ajorlou, Zihui Wu, and Ali Jadbabaie. Demystifying oversmoothing in attention-based graph neural networks. *Advances in Neural Information Processing Systems*, 36:35084–35106, 2023a.
- Xinyi Wu, Zhengdao Chen, William Wei Wang, and Ali Jadbabaie. A non-asymptotic analysis of oversmoothing in graph neural networks. In *The Eleventh International Conference on Learning Representations*, 2023b.
- Zonghan Wu, Shirui Pan, Fengwen Chen, Guodong Long, Chengqi Zhang, and S Yu Philip. A comprehensive survey on graph neural networks. *IEEE transactions on neural networks and learning systems*, 32(1):4–24, 2020.
- Keyulu Xu, Chengtao Li, Yonglong Tian, Tomohiro Sonobe, Ken-ichi Kawarabayashi, and Stefanie Jegelka. Representation learning on graphs with jumping knowledge networks. In *International conference on machine learning*, pp. 5453–5462. PMLR, 2018.
- Zhilin Yang, William Cohen, and Ruslan Salakhutdinov. Revisiting semi-supervised learning with graph embeddings. In *International conference on machine learning*, pp. 40–48. PMLR, 2016.
- Rex Ying, Ruining He, Kaifeng Chen, Pong Eksombatchai, William L Hamilton, and Jure Leskovec. Graph convolutional neural networks for web-scale recommender systems. In *Proceedings of the 24th ACM SIGKDD international conference on knowledge discovery & data mining*, pp. 974–983, 2018.
- Lingxiao Zhao and Leman Akoglu. Pairnorm: Tackling oversmoothing in gnns. In *International Conference on Learning Representations*, 2020.
- Dengyong Zhou, Olivier Bousquet, Thomas Lal, Jason Weston, and Bernhard Schölkopf. Learning with local and global consistency. *Advances in neural information processing systems*, 16, 2003.
- Jie Zhou, Ganqu Cui, Shengding Hu, Zhengyan Zhang, Cheng Yang, Zhiyuan Liu, Lifeng Wang, Changcheng Li, and Maosong Sun. Graph neural networks: A review of methods and applications. *AI Open*, 1:57–81, 2020.
- Kaixiong Zhou, Xiao Huang, Daochen Zha, Rui Chen, Li Li, Soo-Hyun Choi, and Xia Hu. Dirichlet energy constrained learning for deep graph neural networks. *Advances in Neural Information Processing Systems*, 34:21834–21846, 2021.
- Xiaojin Zhu, Zoubin Ghahramani, and John D Lafferty. Semi-supervised learning using gaussian fields and harmonic functions. In *Proceedings of the 20th International conference on Machine learning (ICML-03)*, pp. 912–919, 2003.

A Missing Proofs in Section 3

Proof of Proposition 3.1. We can write matrix \mathbf{H} as:

$$\mathbf{H} = \sum_{i=1}^n \mathbf{H} \mathbf{e}_i \mathbf{e}_i^\top,$$

where each \mathbf{e}_i is the eigenvector of \mathbf{G} associated with eigenvalue λ_i . This indicates that

$$\begin{aligned} \mathbf{H} \tilde{\Delta} &= \mathbf{H}(\mathbf{I} - \mathbf{G}) \\ &= \sum_{i=1}^n (\mathbf{H} \mathbf{e}_i \mathbf{e}_i^\top - \mathbf{H} \mathbf{e}_i \mathbf{e}_i^\top \mathbf{G}) \\ &= \sum_{i=1}^n (1 - \lambda_i) \mathbf{H} \mathbf{e}_i \mathbf{e}_i^\top \\ &= \sum_{i=m+1}^n (1 - \lambda_i) \mathbf{H} \mathbf{e}_i \mathbf{e}_i^\top. \end{aligned}$$

Then using the fact that $1 - \lambda_i \geq 0$ for each i , we obtain

$$\begin{aligned} \|\mathbf{H}\|_E^2 &= \text{Trace}(\mathbf{H} \tilde{\Delta} \mathbf{H}^\top) \\ &= \text{Trace}\left(\sum_{i=m+1}^n (1 - \lambda_i) \mathbf{H} \mathbf{e}_i \mathbf{e}_i^\top \left(\sum_{j=1}^n \mathbf{H} \mathbf{e}_j \mathbf{e}_j^\top\right)^\top\right) \\ &= \text{Trace}\left(\sum_{i=m+1}^n \sum_{j=1}^n (1 - \lambda_i) \mathbf{H} \mathbf{e}_i \mathbf{e}_i^\top \mathbf{e}_j \mathbf{e}_j^\top \mathbf{H}^\top\right) \\ &= \text{Trace}\left(\sum_{i=m+1}^n (1 - \lambda_i) \mathbf{H} \mathbf{e}_i \mathbf{e}_i^\top \mathbf{e}_i \mathbf{e}_i^\top \mathbf{H}^\top\right) \\ &= \text{Trace}\left(\sum_{i=m+1}^n \sqrt{1 - \lambda_i} \mathbf{H} \mathbf{e}_i \mathbf{e}_i^\top \left(\sum_{j=m+1}^n \sqrt{1 - \lambda_j} \mathbf{H} \mathbf{e}_j \mathbf{e}_j^\top\right)^\top\right) \\ &= \left\| \sum_{i=m+1}^n \sqrt{1 - \lambda_i} \mathbf{H} \mathbf{e}_i \mathbf{e}_i^\top \right\|_F^2. \end{aligned}$$

That is,

$$\|\mathbf{H}\|_E = \left\| \sum_{i=m+1}^n \sqrt{1 - \lambda_i} \mathbf{H} \mathbf{e}_i \mathbf{e}_i^\top \right\|_F.$$

On the other hand, equation 3 implies

$$\|\mathbf{H}\|_{\mathcal{M}^\perp} = \|\mathbf{H}_{\mathcal{M}^\perp}\|_F = \left\| \sum_{i=m+1}^n \mathbf{H} \mathbf{e}_i \mathbf{e}_i^\top \right\|_F.$$

We first show that both $\|\mathbf{H}\|_{\mathcal{M}^\perp}$ and $\|\mathbf{H}\|_E$ are seminorms. Since $\|c\mathbf{H}\|_F = |c| \cdot \|\mathbf{H}\|_F, \forall c \in \mathbb{R}$, we have $\|c\mathbf{H}\|_{\mathcal{M}^\perp} = |c| \cdot \|\mathbf{H}\|_{\mathcal{M}^\perp}$ and $\|c\mathbf{H}\|_E = |c| \cdot \|\mathbf{H}\|_E$. Moreover, for any two matrices \mathbf{H}^1 and \mathbf{H}^2 s.t.

$\mathbf{H} = \mathbf{H}^1 + \mathbf{H}^2$, we have

$$\begin{aligned} \sum_{i=m+1}^n \mathbf{H}^1 \mathbf{e}_i \mathbf{e}_i^\top + \sum_{i=m+1}^n \mathbf{H}^2 \mathbf{e}_i \mathbf{e}_i^\top &= \sum_{i=m+1}^n \mathbf{H} \mathbf{e}_i \mathbf{e}_i^\top, \\ \sum_{i=m+1}^n \sqrt{1-\lambda_i} \mathbf{H}^1 \mathbf{e}_i \mathbf{e}_i^\top + \sum_{i=m+1}^n \sqrt{1-\lambda_i} \mathbf{H}^2 \mathbf{e}_i \mathbf{e}_i^\top & \\ &= \sum_{i=m+1}^n \sqrt{1-\lambda_i} \mathbf{H} \mathbf{e}_i \mathbf{e}_i^\top. \end{aligned}$$

Then the triangle inequality of $\|\cdot\|_F$ implies that of $\|\mathbf{H}\|_{\mathcal{M}^\perp}$ and $\|\mathbf{H}\|_E$, respectively.

Now since $0 < 1 - \lambda_{m+1} \leq 1 - \lambda_i \leq 2$ for any $i = m+1, \dots, n$, we may take $\alpha = \sqrt{1 - \lambda_{m+1}}$ and $\beta = \sqrt{2}$. Then

$$\begin{aligned} \alpha \|\mathbf{H}\|_{\mathcal{M}^\perp} &= \left\| \alpha \sum_{i=m+1}^n \mathbf{H} \mathbf{e}_i \mathbf{e}_i^\top \right\|_F \\ &\leq \left\| \sum_{i=m+1}^n \sqrt{1-\lambda_i} \mathbf{H} \mathbf{e}_i \mathbf{e}_i^\top \right\|_F \\ &\leq \left\| \beta \sum_{i=m+1}^n \mathbf{H} \mathbf{e}_i \mathbf{e}_i^\top \right\|_F \\ &= \beta \|\mathbf{H}\|_{\mathcal{M}^\perp}. \end{aligned}$$

The result thus follows from $\|\mathbf{H}\|_E = \left\| \sum_{i=m+1}^n \sqrt{1-\lambda_i} \mathbf{H} \mathbf{e}_i \mathbf{e}_i^\top \right\|_F$. □

A.1 ReLU

Lemma A.1. *Let $\mathbf{Z} \in \mathbb{R}^{d \times n}$, and let $\mathbf{Z}^+ = \max(\mathbf{Z}, 0)$ and $\mathbf{Z}^- = \max(-\mathbf{Z}, 0)$ be the positive and negative parts of \mathbf{Z} . Then (1) $\mathbf{Z}^+, \mathbf{Z}^-$ are (component-wise) nonnegative and $\mathbf{Z} = \mathbf{Z}^+ - \mathbf{Z}^-$ and (2) $\langle \mathbf{Z}^+, \mathbf{Z}^- \rangle_F = 0$.*

Proof of Lemma A.1. Notice that for any $a \in \mathbb{R}$, we have

$$\max(a, 0) = \begin{cases} a & \text{if } a \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

and

$$\max(-a, 0) = \begin{cases} 0 & \text{if } a \geq 0 \\ -a & \text{otherwise} \end{cases}$$

This implies that $a = \max(a, 0) - \max(-a, 0)$ and $\max(a, 0) \cdot \max(-a, 0) = 0$.

Let Z_{ij} be the $(i, j)^{th}$ entry of \mathbf{Z} . Then $\mathbf{Z} = \mathbf{Z}^+ - \mathbf{Z}^-$ follows from $Z_{ij} = \max(Z_{ij}, 0) - \max(-Z_{ij}, 0)$. Also,

$$\begin{aligned} \langle \mathbf{Z}^+, \mathbf{Z}^- \rangle_F &= \text{Trace}((\mathbf{Z}^+)^\top \mathbf{Z}^-) \\ &= \sum_{i=1}^d \sum_{j=1}^n \max(Z_{ij}, 0) \max(-Z_{ij}, 0) = 0. \end{aligned}$$

□

Before proving Proposition 3.2, we notice the following relation between \mathbf{Z} and \mathbf{H} .

Lemma A.2. Given $\mathbf{Z} \in \mathbb{R}^{d \times n}$, let $\mathbf{H} = \sigma(\mathbf{Z})$ with σ being ReLU, then \mathbf{H} lies on the high-dimensional sphere, in $\|\cdot\|_F$ norm, that is centered at $\mathbf{Z}/2$ and with radius $\|\mathbf{Z}/2\|_F$. That is, \mathbf{H} and \mathbf{Z} satisfy the following equation:

$$\left\| \mathbf{H} - \frac{\mathbf{Z}}{2} \right\|_F^2 = \left\| \frac{\mathbf{Z}}{2} \right\|_F^2. \quad (8)$$

Proof of Lemma A.2. We observe that $\mathbf{H} = \sigma(\mathbf{Z}) = \max(\mathbf{Z}, 0) = \mathbf{Z}^+$ is the positive part of \mathbf{Z} . Then we have

$$\langle \mathbf{H}, \mathbf{Z} \rangle_F = \langle \mathbf{H}, \mathbf{Z}^+ - \mathbf{Z}^- \rangle_F = \langle \mathbf{H}, \mathbf{Z}^+ \rangle_F - \langle \mathbf{H}, \mathbf{Z}^- \rangle_F = \langle \mathbf{H}, \mathbf{H} \rangle_F,$$

where we have used $\mathbf{Z} = \mathbf{Z}^+ - \mathbf{Z}^-$ and $\langle \mathbf{H}, \mathbf{Z}^- \rangle_F = \langle \mathbf{Z}^+, \mathbf{Z}^- \rangle_F = 0$ from Lemma A.1.

Therefore, one can deduce the desired result as follows:

$$\begin{aligned} \langle \mathbf{H}, \mathbf{H} \rangle_F - \langle \mathbf{H}, \mathbf{Z} \rangle_F = 0 &\Rightarrow \|\mathbf{H}\|_F^2 - 2 \left\langle \mathbf{H}, \frac{\mathbf{Z}}{2} \right\rangle_F + \left\| \frac{\mathbf{Z}}{2} \right\|_F^2 = \left\| \frac{\mathbf{Z}}{2} \right\|_F^2 \\ &\Rightarrow \left\| \mathbf{H} - \frac{\mathbf{Z}}{2} \right\|_F^2 = \left\| \frac{\mathbf{Z}}{2} \right\|_F^2. \end{aligned}$$

□

Applying $\|\mathbf{H}\|_F^2 = \|\mathbf{H}_{\mathcal{M}} + \mathbf{H}_{\mathcal{M}^\perp}\|_F^2 = \|\mathbf{H}_{\mathcal{M}}\|_F^2 + \|\mathbf{H}_{\mathcal{M}^\perp}\|_F^2$, to both $\frac{\mathbf{Z}}{2}$ and $\mathbf{H} - \frac{\mathbf{Z}}{2}$, we obtain

$$\left\| \frac{\mathbf{Z}}{2} \right\|_F^2 = \left\| \frac{\mathbf{Z}_{\mathcal{M}^\perp}}{2} \right\|_F^2 + \left\| \frac{\mathbf{Z}_{\mathcal{M}}}{2} \right\|_F^2,$$

and

$$\left\| \mathbf{H} - \frac{\mathbf{Z}}{2} \right\|_F^2 = \left\| \mathbf{H}_{\mathcal{M}^\perp} - \frac{\mathbf{Z}_{\mathcal{M}^\perp}}{2} \right\|_F^2 + \left\| \mathbf{H}_{\mathcal{M}} - \frac{\mathbf{Z}_{\mathcal{M}}}{2} \right\|_F^2.$$

Then equation 8 becomes

$$\left\| \frac{\mathbf{Z}_{\mathcal{M}^\perp}}{2} \right\|_F^2 - \left\| \mathbf{H}_{\mathcal{M}^\perp} - \frac{\mathbf{Z}_{\mathcal{M}^\perp}}{2} \right\|_F^2 = \left\| \mathbf{H}_{\mathcal{M}} - \frac{\mathbf{Z}_{\mathcal{M}}}{2} \right\|_F^2 - \left\| \frac{\mathbf{Z}_{\mathcal{M}}}{2} \right\|_F^2 \quad (9)$$

By direct calculation, we have

$$\begin{aligned} \left\| \mathbf{H}_{\mathcal{M}} - \frac{\mathbf{Z}_{\mathcal{M}}}{2} \right\|_F^2 - \left\| \frac{\mathbf{Z}_{\mathcal{M}}}{2} \right\|_F^2 &= \langle \mathbf{H}_{\mathcal{M}}, \mathbf{H}_{\mathcal{M}} \rangle_F - 2 \left\langle \mathbf{H}_{\mathcal{M}}, \frac{\mathbf{Z}_{\mathcal{M}}}{2} \right\rangle_F \\ &= \langle \mathbf{H}_{\mathcal{M}}, \mathbf{H}_{\mathcal{M}} - \mathbf{Z}_{\mathcal{M}} \rangle_F. \end{aligned} \quad (10)$$

Combining equation 9 and equation 10, we obtain the following result:

Lemma A.3. For any $\mathbf{Z} = \mathbf{Z}_{\mathcal{M}} + \mathbf{Z}_{\mathcal{M}^\perp}$, let

$$\mathbf{H} = \sigma(\mathbf{Z}) = \mathbf{H}_{\mathcal{M}} + \mathbf{H}_{\mathcal{M}^\perp},$$

then we have

$$\left\| \frac{\mathbf{Z}_{\mathcal{M}^\perp}}{2} \right\|_F^2 - \left\| \mathbf{H}_{\mathcal{M}^\perp} - \frac{\mathbf{Z}_{\mathcal{M}^\perp}}{2} \right\|_F^2 = \langle \mathbf{Z}_{\mathcal{M}}^+, \mathbf{Z}_{\mathcal{M}}^- \rangle_F.$$

where $\mathbf{Z}_{\mathcal{M}}^+ = \sum_{i=1}^m \mathbf{Z}^+ \mathbf{e}_i \mathbf{e}_i^\top$, $\mathbf{Z}_{\mathcal{M}}^- = \sum_{i=1}^m \mathbf{Z}^- \mathbf{e}_i \mathbf{e}_i^\top$.

Proof of Lemma A.3. Recall that $\mathbf{H} = \sigma(\mathbf{Z}) = \max(\mathbf{Z}, 0) = \mathbf{Z}^+$. Also, $\mathbf{Z} = \mathbf{Z}^+ - \mathbf{Z}^-$ implies $\mathbf{Z}_{\mathcal{M}} = \mathbf{Z}_{\mathcal{M}}^+ - \mathbf{Z}_{\mathcal{M}}^- = \mathbf{H}_{\mathcal{M}}^+ - \mathbf{Z}_{\mathcal{M}}^-$. Therefore, we see that

$$\langle \mathbf{H}_{\mathcal{M}}, \mathbf{H}_{\mathcal{M}} - \mathbf{Z}_{\mathcal{M}} \rangle_F = \langle \mathbf{Z}_{\mathcal{M}}^+, \mathbf{Z}_{\mathcal{M}}^- \rangle_F.$$

□

By using the fact that $\langle \mathbf{Z}_{\mathcal{M}}^+, \mathbf{Z}_{\mathcal{M}}^- \rangle_F \geq 0$ in Lemma A.3, we reveal a geometric relation between \mathbf{Z} and \mathbf{H} mentioned in Proposition 3.2.

Proof of Proposition 3.2. Since $\mathbf{Z}^+, \mathbf{Z}^- \geq 0$ are nonnegative and all the eigenvectors \mathbf{e}_i are also nonnegative, we see that $\mathbf{Z}_{\mathcal{M}}^+ = \sum_{i=1}^m \mathbf{Z}^+ \mathbf{e}_i \mathbf{e}_i^\top$ and $\mathbf{Z}_{\mathcal{M}}^- = \sum_{i=1}^m \mathbf{Z}^- \mathbf{e}_i \mathbf{e}_i^\top$ are nonnegative. This indicates that

$$\langle \mathbf{Z}_{\mathcal{M}}^+, \mathbf{Z}_{\mathcal{M}}^- \rangle_F = \text{Trace}(\mathbf{Z}_{\mathcal{M}}^+ (\mathbf{Z}_{\mathcal{M}}^-)^\top) \geq 0.$$

Then according to Lemma A.3, we obtain

$$\left\| \frac{\mathbf{Z}_{\mathcal{M}^\perp}}{2} \right\|_F^2 - \left\| \mathbf{H}_{\mathcal{M}^\perp} - \frac{\mathbf{Z}_{\mathcal{M}^\perp}}{2} \right\|_F^2 = \langle \mathbf{Z}_{\mathcal{M}}^+, \mathbf{Z}_{\mathcal{M}}^- \rangle_F \geq 0.$$

So we have

$$\begin{aligned} \left\| \mathbf{H}_{\mathcal{M}^\perp} - \frac{\mathbf{Z}_{\mathcal{M}^\perp}}{2} \right\|_F &= \sqrt{\left\| \frac{\mathbf{Z}_{\mathcal{M}^\perp}}{2} \right\|_F^2 - \langle \mathbf{Z}_{\mathcal{M}}^+, \mathbf{Z}_{\mathcal{M}}^- \rangle_F} \\ &= \sqrt{\left\| \frac{\mathbf{Z}_{\mathcal{M}^\perp}}{2} \right\|_F^2 - \langle \mathbf{H}_{\mathcal{M}}, \mathbf{H}_{\mathcal{M}} - \mathbf{Z}_{\mathcal{M}} \rangle_F}, \end{aligned}$$

which shows that $\mathbf{H}_{\mathcal{M}^\perp}$ lies on the high-dimensional sphere that we have claimed. Furthermore, we conclude that

$$0 \leq \left\| \mathbf{H}_{\mathcal{M}^\perp} - \frac{\mathbf{Z}_{\mathcal{M}^\perp}}{2} \right\|_F \leq \left\| \frac{\mathbf{Z}_{\mathcal{M}^\perp}}{2} \right\|_F. \quad (11)$$

This demonstrates that $\mathbf{H}_{\mathcal{M}^\perp}$ lies on the high-dimensional sphere we have stated.

Since the sphere $\left\| \mathbf{H}_{\mathcal{M}^\perp} - \frac{\mathbf{Z}_{\mathcal{M}^\perp}}{2} \right\|_F^2 = \left\| \frac{\mathbf{Z}_{\mathcal{M}^\perp}}{2} \right\|_F^2$ passes through the origin, the distance of any $\mathbf{H}_{\mathcal{M}^\perp}$ to the origin must be no greater than the diameter of this sphere, i.e., $\|\mathbf{H}_{\mathcal{M}^\perp}\|_F \leq \|\mathbf{Z}_{\mathcal{M}^\perp}\|_F$. Also, this can be derived from the following inequality:

$$\|\mathbf{H}_{\mathcal{M}^\perp}\|_F - \left\| \frac{\mathbf{Z}_{\mathcal{M}^\perp}}{2} \right\|_F \leq \left\| \mathbf{H}_{\mathcal{M}^\perp} - \frac{\mathbf{Z}_{\mathcal{M}^\perp}}{2} \right\|_F \leq \left\| \frac{\mathbf{Z}_{\mathcal{M}^\perp}}{2} \right\|_F.$$

We see that the maximal smoothness $\|\mathbf{H}_{\mathcal{M}^\perp}\|_F = \|\mathbf{Z}_{\mathcal{M}^\perp}\|_F$ is attained when $\mathbf{H}_{\mathcal{M}^\perp} = \mathbf{Z}_{\mathcal{M}^\perp}$, the intersection of the surface and the line passing through the center and the origin.

After all, we complete the proof by using the fact that $\|\mathbf{Z}_{\mathcal{M}^\perp}\|_F = \|\mathbf{Z}\|_{\mathcal{M}^\perp}$ for any matrix \mathbf{Z} , which implies $\|\mathbf{H}\|_{\mathcal{M}^\perp} = \|\mathbf{H}_{\mathcal{M}^\perp}\|_F \leq \|\mathbf{Z}_{\mathcal{M}^\perp}\|_F = \|\mathbf{Z}\|_{\mathcal{M}^\perp}$. \square

A.2 Leaky ReLU

For the leaky ReLU activation, we have the following Lemma:

Lemma A.4. *If $\mathbf{H} = \sigma_a(\mathbf{Z})$ with σ_a being the leaky ReLU activation function, then \mathbf{H} lies on the high-dimensional sphere centered at $(1+a)\mathbf{Z}/2$ with radius $\|(1-a)\mathbf{Z}/2\|_F$.*

Proof of Lemma A.4. Notice that

$$\mathbf{H} = \sigma_a(\mathbf{Z}) = \mathbf{Z}^+ - a\mathbf{Z}^-.$$

Then $\mathbf{H} - \mathbf{Z} = (1-a)\mathbf{Z}^-$ and $\mathbf{H} - a\mathbf{Z} = (1-a)\mathbf{Z}^+$. Using $\langle \mathbf{Z}^-, \mathbf{Z}^+ \rangle_F = 0$, we have

$$\begin{aligned} \langle \mathbf{H} - \mathbf{Z}, \mathbf{H} - a\mathbf{Z} \rangle_F &= 0 \\ \Rightarrow \|\mathbf{H}\|_F^2 - 2 \left\langle \mathbf{H}, \frac{(1+a)\mathbf{Z}}{2} \right\rangle_F + a\|\mathbf{Z}\|_F^2 &= 0 \\ \Rightarrow \|\mathbf{H}\|_F^2 - 2 \left\langle \mathbf{H}, \frac{(1+a)\mathbf{Z}}{2} \right\rangle_F &= -a\|\mathbf{Z}\|_F^2 \\ \Rightarrow \left\| \mathbf{H} - \frac{(1+a)\mathbf{Z}}{2} \right\|_F^2 &= \left\| \frac{(1+a)\mathbf{Z}}{2} \right\|_F^2 - a\|\mathbf{Z}\|_F^2 = \left\| \frac{(1-a)\mathbf{Z}}{2} \right\|_F^2. \end{aligned}$$

\square

Moreover, we notice that

Lemma A.5. For any $\mathbf{Z} = \mathbf{Z}_{\mathcal{M}} + \mathbf{Z}_{\mathcal{M}^\perp}$, let $\mathbf{H} = \sigma_a(\mathbf{Z}) = \mathbf{H}_{\mathcal{M}} + \mathbf{H}_{\mathcal{M}^\perp}$, then

$$\begin{aligned} & \left\| \frac{(1-a)}{2} \mathbf{Z}_{\mathcal{M}^\perp} \right\|_F^2 - \left\| \mathbf{H}_{\mathcal{M}^\perp} - \frac{(1+a)}{2} \mathbf{Z}_{\mathcal{M}^\perp} \right\|_F^2 \\ &= (1-a)^2 \langle \mathbf{Z}_{\mathcal{M}^\perp}^+, \mathbf{Z}_{\mathcal{M}^\perp}^- \rangle_F. \end{aligned}$$

Proof of Lemma A.5. Similar to the proof of Lemma A.3, the orthogonal decomposition implies that

$$\begin{aligned} & \left\| \frac{(1-a)}{2} \mathbf{Z}_{\mathcal{M}^\perp} \right\|_F^2 - \left\| \mathbf{H}_{\mathcal{M}^\perp} - \frac{(1+a)}{2} \mathbf{Z}_{\mathcal{M}^\perp} \right\|_F^2 \\ &= \left\| \mathbf{H}_{\mathcal{M}} - \frac{(1+a)}{2} \mathbf{Z}_{\mathcal{M}} \right\|_F^2 - \left\| \frac{(1-a)}{2} \mathbf{Z}_{\mathcal{M}} \right\|_F^2 \\ &= \langle \mathbf{H}_{\mathcal{M}} - \mathbf{Z}_{\mathcal{M}}, \mathbf{H}_{\mathcal{M}} - a \mathbf{Z}_{\mathcal{M}} \rangle_F \\ &= \langle (1-a) \mathbf{Z}_{\mathcal{M}}^-, (1-a) \mathbf{Z}_{\mathcal{M}}^+ \rangle_F \\ &= (1-a)^2 \langle \mathbf{Z}_{\mathcal{M}}^-, \mathbf{Z}_{\mathcal{M}}^+ \rangle_F. \end{aligned}$$

□

Proof of Proposition 3.3. Similar to the proof of Proposition 3.2, we apply $\langle \mathbf{Z}_{\mathcal{M}}^-, \mathbf{Z}_{\mathcal{M}}^+ \rangle_F \geq 0$ to Lemma A.5 and hence obtain the geometric condition as follows:

$$\left\| \mathbf{H}_{\mathcal{M}^\perp} - \frac{(1+a)}{2} \mathbf{Z}_{\mathcal{M}^\perp} \right\|_F = \sqrt{\left\| \frac{(1-a)}{2} \mathbf{Z}_{\mathcal{M}^\perp} \right\|_F^2 - \langle \mathbf{H}_{\mathcal{M}} - \mathbf{Z}_{\mathcal{M}}, \mathbf{H}_{\mathcal{M}} - a \mathbf{Z}_{\mathcal{M}} \rangle_F}.$$

Then we have the following inequality:

$$0 \leq \left\| \mathbf{H}_{\mathcal{M}^\perp} - \frac{(1+a)}{2} \mathbf{Z}_{\mathcal{M}^\perp} \right\|_F \leq \left\| \frac{(1-a)}{2} \mathbf{Z}_{\mathcal{M}^\perp} \right\|_F.$$

Moreover, we deduce that

$$\begin{aligned} \left| \left\| \mathbf{H}_{\mathcal{M}^\perp} \right\|_F - \left\| \frac{(1+a)}{2} \mathbf{Z}_{\mathcal{M}^\perp} \right\|_F \right| &\leq \left\| \mathbf{H}_{\mathcal{M}^\perp} - \frac{(1+a)}{2} \mathbf{Z}_{\mathcal{M}^\perp} \right\|_F \\ &\leq \left\| \frac{(1-a)}{2} \mathbf{Z}_{\mathcal{M}^\perp} \right\|_F. \end{aligned}$$

and hence

$$\begin{aligned} - \left\| \frac{(1-a)}{2} \mathbf{Z}_{\mathcal{M}^\perp} \right\|_F &\leq \left\| \mathbf{H}_{\mathcal{M}^\perp} \right\|_F - \left\| \frac{(1+a)}{2} \mathbf{Z}_{\mathcal{M}^\perp} \right\|_F \\ &\leq \left\| \frac{(1-a)}{2} \mathbf{Z}_{\mathcal{M}^\perp} \right\|_F. \end{aligned}$$

Therefore, we obtain $a \left\| \mathbf{Z}_{\mathcal{M}^\perp} \right\|_F \leq \left\| \mathbf{H}_{\mathcal{M}^\perp} \right\|_F \leq \left\| \mathbf{Z}_{\mathcal{M}^\perp} \right\|_F$. (Remark that $\mathbf{H}_{\mathcal{M}^\perp}$ achieves its maximal norm when it is equal to $\mathbf{Z}_{\mathcal{M}^\perp}$, the intersection of the surface and the line passing through the center and the origin.)

By using the fact that $\left\| \mathbf{Z}_{\mathcal{M}^\perp} \right\|_F = \left\| \mathbf{Z} \right\|_{\mathcal{M}^\perp}$ for any matrix \mathbf{Z} , we conclude that $a \left\| \mathbf{Z} \right\|_{\mathcal{M}^\perp} \leq \left\| \mathbf{H} \right\|_{\mathcal{M}^\perp} \leq \left\| \mathbf{Z} \right\|_{\mathcal{M}^\perp}$. □

B Proofs in Section 4

Throughout this section, we assume that $\mathbf{z}_{\mathcal{M}^\perp} \neq \mathbf{0}$.

Proof of Proposition 4.3. Recall that $\mathbf{e} = \tilde{\mathbf{D}}^{\frac{1}{2}} \mathbf{u}_n / c$ has only positive entries where $\tilde{\mathbf{D}}$ is the augmented degree matrix and $\mathbf{u}_n = [1, \dots, 1]^\top \in \mathbb{R}^n$ and $c = \left\| \tilde{\mathbf{D}}^{\frac{1}{2}} \mathbf{u}_n \right\|$. Let d_i be the i^{th} diagonal entry of $\tilde{\mathbf{D}}$. Then we have $\mathbf{e} = [\sqrt{d_1}/c, \sqrt{d_2}/c, \dots, \sqrt{d_n}/c]^\top$ and $c = \sqrt{\sum_{i=1}^n d_i}$.

Note that $\mathbf{z}(\alpha) = \mathbf{z} - \alpha \mathbf{e} = \mathbf{z} - \frac{\alpha}{c} \tilde{\mathbf{D}}^{\frac{1}{2}} \mathbf{u}_n = \tilde{\mathbf{D}}^{\frac{1}{2}} (\tilde{\mathbf{D}}^{-\frac{1}{2}} \mathbf{z} - \frac{\alpha}{c} \mathbf{u}_n) = \tilde{\mathbf{D}}^{\frac{1}{2}} (\mathbf{x} - \frac{\alpha}{c} \mathbf{u}_n)$, where we assume $\mathbf{x} := \tilde{\mathbf{D}}^{-\frac{1}{2}} \mathbf{z}$. Then we observe that when σ is the ReLU activation function,

$$\mathbf{h}(\alpha) = \sigma(\mathbf{z}(\alpha)) = \sigma\left(\tilde{\mathbf{D}}^{\frac{1}{2}} \left(\mathbf{x} - \frac{\alpha}{c} \mathbf{u}_n\right)\right) = \tilde{\mathbf{D}}^{\frac{1}{2}} \sigma\left(\mathbf{x} - \frac{\alpha}{c} \mathbf{u}_n\right),$$

and hence

$$\langle \mathbf{h}(\alpha), \mathbf{e} \rangle = \left\langle \tilde{\mathbf{D}}^{\frac{1}{2}} \sigma\left(\mathbf{x} - \frac{\alpha}{c} \mathbf{u}_n\right), \mathbf{e} \right\rangle = \left\langle \sigma\left(\mathbf{x} - \frac{\alpha}{c} \mathbf{u}_n\right), \tilde{\mathbf{D}} \mathbf{u}_n \right\rangle.$$

We may now assume $\mathbf{x} = [x_1, \dots, x_n]^\top$ is well-ordered s.t. $x_1 \geq x_2 \geq \dots \geq x_n$. Indeed, there is a collection of indices $\{k_1, \dots, k_l\}$ such that

$$\begin{aligned} x_1 &= \dots, x_{k_1} \text{ and } x_{k_1} > x_{k_1+1}, \\ x_{k_{j-1}+1} &= \dots = x_{k_j} \text{ and } x_{k_j} > x_{k_j+1} \text{ for any } j = 2, \dots, l-1, \\ x_{k_{l-1}+1} &= \dots = x_{k_l} \text{ and } k_l = n. \end{aligned}$$

That is, $x_1 = x_2 = \dots = x_{k_1} > x_{k_1+1} = \dots = x_{k_2} > x_{k_2+1} = \dots = x_{k_3} > x_{k_3+1} \dots$

We first restrict the domain of α s.t. $\mathbf{h}(\alpha) \neq 0$. Note that we have

$$\begin{aligned} \mathbf{h}(\alpha) = 0 &\Leftrightarrow \sigma\left(\mathbf{x} - \frac{\alpha}{c} \mathbf{u}_n\right) = 0 \\ &\Leftrightarrow x_i - \frac{\alpha}{c} \leq 0 \text{ for } i = 1, \dots, n \\ &\Leftrightarrow x_1 - \frac{\alpha}{c} \leq 0 \\ &\Leftrightarrow \alpha \geq cx_1. \end{aligned}$$

So we will study the smoothness $s(\mathbf{h}(\alpha))$ when $\alpha < cx_1$.

Let $\epsilon > 0$ and consider $\alpha = c(x_1 - \epsilon)$. When $\epsilon \leq x_1 - x_{k_1+1} = x_1 - x_{k_2}$, we see that

$$\mathbf{x} - \frac{\alpha}{c} \mathbf{u}_n = [\epsilon, \dots, \epsilon, \epsilon - (x_1 - x_{k_1+1}), \dots, \epsilon - (x_1 - x_n)]^\top,$$

where only the first k_1 entries are positive since $x_1 - x_i \geq \epsilon$ for any $i \geq k_1 + 1$. Therefore,

$$\begin{aligned} \mathbf{h}(\alpha) &= \tilde{\mathbf{D}}^{\frac{1}{2}} \sigma\left(\mathbf{x} - \frac{\alpha}{c} \mathbf{u}_n\right) \\ &= \tilde{\mathbf{D}}^{\frac{1}{2}} [\epsilon, \dots, \epsilon, 0, \dots, 0]^\top \\ &= [\epsilon \sqrt{d_1}, \dots, \epsilon \sqrt{d_{k_1}}, 0, \dots, 0]^\top. \end{aligned}$$

Hence we have $\|\mathbf{h}(\alpha)\| = \epsilon \sqrt{\sum_{i=1}^{k_1} d_i}$. Also, we have

$$\begin{aligned} \|\mathbf{h}(\alpha)\|_{\mathcal{M}} &= |\langle \mathbf{h}(\alpha), \mathbf{e} \rangle| \\ &= [\epsilon \sqrt{d_1}, \dots, \epsilon \sqrt{d_{k_1}}, 0, \dots, 0]^\top [\sqrt{d_1}/c, \sqrt{d_2}/c, \dots, \sqrt{d_n}/c] \\ &= \frac{\epsilon}{c} \sum_{i=1}^{k_1} d_i. \end{aligned}$$

Then we obtain the smoothness $s(\mathbf{h}(\alpha))$ as follows

$$s(\mathbf{h}(\alpha)) = \frac{\|\mathbf{h}(\alpha)\|_{\mathcal{M}}}{\|\mathbf{h}(\alpha)\|} = \frac{\frac{\epsilon}{c} \sum_{i=1}^{k_1} d_i}{\epsilon \sqrt{\sum_{i=1}^{k_1} d_i}} = \frac{\sqrt{\sum_{i=1}^{k_1} d_i}}{c} = \frac{K_1}{c} < 1,$$

where $K_1 := \sqrt{\sum_{i=1}^{k_1} d_i}$. Similarly, we may denote $\sqrt{\sum_{i=k_{j-1}+1}^{k_j} d_i}$ by K_j for $j = 2, \dots, l$.

Now we are going to show that the smoothness $s(\mathbf{h}(\alpha))$ is increasing as α gets smaller whenever $\alpha < c x_1$, implying $\frac{K_1}{c}$ is the minimum of the smoothness $s(\mathbf{h}(\alpha))$. Remember that we are considering $\alpha = c(x_1 - \epsilon)$ and we have studied the case when $0 < \epsilon \leq x_1 - x_{k_1+1} = x_1 - x_{k_2}$.

Let $\delta_j := x_1 - x_{k_j}$ for $1 \leq j \leq l$. Clearly, we have $\delta_1 = 0$ and $\delta_j < \delta_{j+1}$ for $1 \leq j \leq l-1$. Fix a $j' \in \{2, \dots, l-1\}$, we see that when $\delta_{j'} < \epsilon \leq x_1 - x_{k_{j'}+1}$,

$$\begin{aligned} \mathbf{x} - \frac{\alpha}{c} \mathbf{u}_n &= \left[\epsilon - \delta_1, \dots, \epsilon - \delta_1, \epsilon - \delta_2, \dots, \epsilon - \delta_2, \epsilon - \delta_3, \right. \\ &\quad \left. \dots, \epsilon - \delta_{j'}, \epsilon - (x_1 - x_{k_{j'}+1}), \dots, \epsilon - (x_1 - x_n) \right]^\top, \end{aligned}$$

where we have $\epsilon - \delta_j > 0$ for $2 \leq j \leq j'$ and $\epsilon - (x_1 - x_i) \leq 0$ for any $i \geq k_{j'} + 1$. Consequently,

$$\begin{aligned} \mathbf{h}(\alpha) &= \tilde{\mathbf{D}}^{\frac{1}{2}} \sigma \left(\mathbf{x} - \frac{\alpha}{c} \mathbf{u}_n \right) \\ &= [(\epsilon - \delta_1) \sqrt{d_1}, \dots, (\epsilon - \delta_1) \sqrt{d_{k_1}}, (\epsilon - \delta_2) \sqrt{d_{k_1+1}}, \dots, \\ &\quad (\epsilon - \delta_2) \sqrt{d_{k_2}}, (\epsilon - \delta_3) \sqrt{d_{k_2+1}}, \dots, (\epsilon - \delta_{j'}) \sqrt{d_{k_{j'}}}, 0, \dots, 0]^\top. \end{aligned}$$

Then we can compute

$$\|\mathbf{h}(\alpha)\| = \sqrt{\sum_{j=1}^{j'} \sum_{i=k_{j-1}+1}^{k_j} d_i (\epsilon - \delta_j)^2} = \sqrt{\sum_{j=1}^{j'} K_j^2 (\epsilon - \delta_j)^2},$$

where we set $k_0 := 0$ for simplicity and $K_j = \sqrt{\sum_{i=k_{j-1}+1}^{k_j} d_i}$ for $j = 1, \dots, j'$. Also, we have

$$\begin{aligned} \|\mathbf{h}(\alpha)\|_{\mathcal{M}} &= |\langle \mathbf{h}(\alpha), \mathbf{e} \rangle| = \sum_{j=1}^{j'} \sum_{i=k_{j-1}+1}^{k_j} \frac{d_i (\epsilon - \delta_j)}{c} \\ &= \frac{1}{c} \sum_{j=1}^{j'} K_j^2 (\epsilon - \delta_j). \end{aligned}$$

A careful calculation shows that $\frac{\partial}{\partial \epsilon} s(\mathbf{h}(\alpha)) > 0$ whenever $\delta_{j'} < \epsilon \leq x_1 - x_{k_{j'}+1}$ which implies that $s(\mathbf{h}(\alpha))$ is increasing as ϵ increases. Indeed, we have

$$\begin{aligned} &\frac{\partial}{\partial \epsilon} s(\mathbf{h}(\alpha)) \\ &= \frac{\partial}{\partial \epsilon} \left(\frac{\sum_{j=1}^{j'} K_j^2 (\epsilon - \delta_j)}{c \sqrt{\sum_{j=1}^{j'} K_j^2 (\epsilon - \delta_j)^2}} \right) \\ &= \frac{\left(\frac{\partial}{\partial \epsilon} \sum_{j=1}^{j'} K_j^2 (\epsilon - \delta_j) \right) \sqrt{\sum_{j=1}^{j'} K_j^2 (\epsilon - \delta_j)^2}}{c \sum_{j=1}^{j'} K_j^2 (\epsilon - \delta_j)^2} - \\ &\quad - \frac{\sum_{j=1}^{j'} K_j^2 (\epsilon - \delta_j) \left(\frac{\partial}{\partial \epsilon} \sqrt{\sum_{j=1}^{j'} K_j^2 (\epsilon - \delta_j)^2} \right)}{c \sum_{j=1}^{j'} K_j^2 (\epsilon - \delta_j)^2} \\ &= \frac{\left(\sum_{j=1}^{j'} K_j^2 \right) \sum_{j=1}^{j'} K_j^2 (\epsilon - \delta_j)^2 - \sum_{j=1}^{j'} K_j^2 (\epsilon - \delta_j) \left(\sum_{j=1}^{j'} K_j^2 (\epsilon - \delta_j) \right)}{c \sum_{j=1}^{j'} K_j^2 (\epsilon - \delta_j)^2 \sqrt{\sum_{j=1}^{j'} K_j^2 (\epsilon - \delta_j)^2}}. \end{aligned}$$

Then to show that $\frac{\partial}{\partial \epsilon} s(\mathbf{h}(\alpha)) > 0$, it suffices to show that the numerator is positive, i.e.

$$\left(\sum_{j=1}^{j'} K_j^2 \right) \sum_{j=1}^{j'} K_j^2 (\epsilon - \delta_j)^2 - \left(\sum_{j=1}^{j'} K_j^2 (\epsilon - \delta_j) \right)^2 > 0,$$

since $c \sum_{j=1}^{j'} K_j^2 (\epsilon - \delta_j)^2 \sqrt{\sum_{j=1}^{j'} K_j^2 (\epsilon - \delta_j)^2} > 0$ is always positive. In fact, this follows from the Cauchy inequality $\|\mathbf{v}\| \|\mathbf{u}\| \geq \langle \mathbf{v}, \mathbf{u} \rangle$, where we set

$$\begin{aligned} \mathbf{v} &:= [K_1, K_2, \dots, K_{j'}]^\top, \\ \mathbf{u} &:= [K_1(\epsilon - \delta_1), K_2(\epsilon - \delta_2), \dots, K_{j'}(\epsilon - \delta_{j'})]^\top. \end{aligned}$$

Moreover, equality happens only when \mathbf{v} is parallel to \mathbf{u} . This is, however, impossible since $\epsilon - \delta_j > \epsilon - \delta_{j+1}$ for any $j = 1, \dots, j' - 1$ and each K_j is positive.

So we see that $s(\mathbf{h}(\alpha))$ is increasing as ϵ increases whenever $0 < \epsilon$, and hence the smoothness $s(\mathbf{h}(\alpha))$ is increasing as α decreases whenever $cx_n \leq \alpha < cx_1$.

For the case $j' = l$ where $\delta_l = x_1 - x_n < \epsilon$, we have $x_n - \alpha/c = x_n - (x_1 - \epsilon) = \epsilon - (x_1 - x_n) > 0$, implying $\alpha < cx_n$ and $\mathbf{h}(\alpha) = \mathbf{z}(\alpha)$. We have shown that the smoothness is increasing as α is going far from $\langle \mathbf{z}, \mathbf{e} \rangle$; in particular, when $\alpha < \langle \mathbf{z}, \mathbf{e} \rangle$ and α is decreasing. One can check that

$$\begin{aligned} cx_n &= \frac{\sum_{i=1}^n d_i x_n}{c} = \left\langle x_n \mathbf{u}_n, \frac{\tilde{\mathbf{D}} \mathbf{u}_n}{c} \right\rangle \\ &\leq \left\langle \mathbf{x}, \frac{\tilde{\mathbf{D}} \mathbf{u}_n}{c} \right\rangle \\ &= \left\langle \tilde{\mathbf{D}}^{\frac{1}{2}} \mathbf{x}, \frac{\tilde{\mathbf{D}}^{\frac{1}{2}} \mathbf{u}_n}{c} \right\rangle \\ &= \langle \mathbf{z}, \mathbf{e} \rangle, \end{aligned}$$

which means the smoothness is increasing as α decreases whenever $\alpha < cx_n$.

We conclude that the smoothness increases as α decreases, provided $\alpha < cx_1$. Also, we have $\sup_{\alpha < cx_1} s(\mathbf{h}(\alpha)) = 1$ as the case in the proof of Proposition B.1. One can check that $s(\mathbf{h}(\alpha))$ is a continuous function for $\alpha < cx_1$ and thus it has range $[K_1/c, 1)$ by the mean value theorem.

Finally, we can establish the result: $K_1/c = \sqrt{\frac{\sum_{x_i = \max \mathbf{x}} d_i}{\sum_{j=1}^n d_j}}$ is the minimum of $s(\mathbf{h}(\alpha))$ and 1 is the maximum of $s(\mathbf{h}(\alpha))$ occurring whenever $\alpha \geq cx_1 = \sqrt{\sum_{j=1}^n d_j} \max_i x_i$. Moreover, $s(\mathbf{h}(\alpha))$ has a monotone property when $\alpha < \sqrt{\sum_{j=1}^n d_j} \max_i x_i$ and has range $\left[\sqrt{\frac{\sum_{x_i = \max \mathbf{x}} d_i}{\sum_{j=1}^n d_j}}, 1 \right]$.

It is clear that the assumption on the ordering of the entries of \mathbf{x} will not affect this result. \square

To prove Proposition 4.4, we first prove an analogous result for the identity function, that is, $\mathbf{h} = \sigma(\mathbf{z}) = \mathbf{z}$.

Proposition B.1. *Suppose $\mathbf{z}_{\mathcal{M}^\perp} \neq \mathbf{0}$, then $s(\mathbf{z}(\alpha))$ achieves its minimum 0 if $\alpha = \langle \mathbf{z}, \mathbf{e} \rangle$. Moreover, $\sup_\alpha s(\mathbf{z}(\alpha)) = 1$ where $s(\mathbf{z}(\alpha))$ is close to 1 when α is far away from $\langle \mathbf{z}, \mathbf{e} \rangle$.*

Proof of Proposition B.1. We know that $0 \leq s(\mathbf{z}(\alpha)) \leq 1$ and

$$\begin{aligned} s(\mathbf{z}(\alpha)) &= \sqrt{1 - \frac{\|\mathbf{z}_{\mathcal{M}^\perp}\|^2}{\|\mathbf{z}(\alpha)\|^2}} \\ &= \sqrt{1 - \frac{\|\mathbf{z}_{\mathcal{M}^\perp}\|^2}{\|\mathbf{z}_{\mathcal{M}^\perp}\|^2 + \|\mathbf{z}(\alpha)_{\mathcal{M}}\|^2}} \\ &= \sqrt{1 - \frac{\|\mathbf{z}_{\mathcal{M}^\perp}\|^2}{\|\mathbf{z}_{\mathcal{M}^\perp}\|^2 + \|\mathbf{z}_{\mathcal{M}} - \alpha \mathbf{e}\|^2}}. \end{aligned}$$

Suppose $s(\mathbf{z}(\alpha)) = 1$. Then we have $\frac{\|\mathbf{z}_{\mathcal{M}^\perp}\|^2}{\|\mathbf{z}_{\mathcal{M}^\perp}\|^2 + \|\mathbf{z}_{\mathcal{M}} - \alpha \mathbf{e}\|^2} = 0$ which forces $\|\mathbf{z}_{\mathcal{M}^\perp}\| = 0$. However, this contradicts the hypothesis $\mathbf{z}_{\mathcal{M}^\perp} \neq 0$. So $s(\mathbf{z}(\alpha))$ cannot attain its maximum.

But for any $0 \leq t < 1$, $s(\mathbf{z}(\alpha)) = t$ if and only if

$$\begin{aligned} \sqrt{1 - \frac{\|\mathbf{z}_{\mathcal{M}^\perp}\|^2}{\|\mathbf{z}_{\mathcal{M}^\perp}\|^2 + \|\mathbf{z}_{\mathcal{M}} - \alpha \mathbf{e}\|^2}} &= t \\ \Leftrightarrow \frac{\|\mathbf{z}_{\mathcal{M}^\perp}\|^2}{\|\mathbf{z}_{\mathcal{M}^\perp}\|^2 + \|\mathbf{z}_{\mathcal{M}} - \alpha \mathbf{e}\|^2} &= 1 - t^2 \\ \Leftrightarrow \|\mathbf{z}_{\mathcal{M}^\perp}\|^2 &= (1 - t^2)(\|\mathbf{z}_{\mathcal{M}^\perp}\|^2 + \|\mathbf{z}_{\mathcal{M}} - \alpha \mathbf{e}\|^2) \\ \Leftrightarrow t^2 \|\mathbf{z}_{\mathcal{M}^\perp}\|^2 &= (1 - t^2) \|\mathbf{z}_{\mathcal{M}} - \alpha \mathbf{e}\|^2 \\ \Leftrightarrow \|\mathbf{z}_{\mathcal{M}} - \alpha \mathbf{e}\| &= \sqrt{\frac{t^2}{1 - t^2}} \cdot \|\mathbf{z}_{\mathcal{M}^\perp}\| \end{aligned}$$

This implies that $\sup_\alpha s(\mathbf{z}(\alpha)) = 1$ and $s(\mathbf{z}(\alpha))$ achieves its minimum 0 if and only if $\alpha = \langle \mathbf{z}, \mathbf{e} \rangle$. It is clear that $s(\mathbf{z}(\alpha))$ get closer to 1 when α is going far away from $\langle \mathbf{z}, \mathbf{e} \rangle$. i.e., $|\alpha - \langle \mathbf{z}, \mathbf{e} \rangle| = \|\mathbf{z}_{\mathcal{M}} - \alpha \mathbf{e}\|$ is increasing. \square

Proof of Proposition 4.4. First, we notice that leaky ReLU has the following two properties

1. $\sigma_a(x) > 0$ for $x \gg 0$ and $\sigma_a(x) < 0$ for $x \ll 0$.
2. σ_a is a non-trivial linear map for $x \gg 0$.

We will use Property 1 to show that $\min_\alpha s(\mathbf{h}(\alpha)) = 0$ and Property 2 to show that $\sup_\alpha s(\mathbf{h}(\alpha)) = 1$. Notice that $\sigma_a(x) < 0$ for $x \ll 0$ implies that there exists a sufficient small $\alpha_2 < 0$ s.t. all of the entries of $\mathbf{h}(\alpha_2)$ are negative and hence $|\langle \mathbf{h}(\alpha_2), \mathbf{e} \rangle| < 0$. Similarly, $\sigma_a(x) > 0$ for $x \gg 0$ implies that there exists a sufficient large $\alpha_1 > 0$ s.t. all of the entries of $\mathbf{h}(\alpha_1)$ are positive and hence $|\langle \mathbf{h}(\alpha_1), \mathbf{e} \rangle| > 0$. Since $|\langle \mathbf{h}(\alpha), \mathbf{e} \rangle|$ is a continuous function of α on $[\alpha_1, \alpha_2]$, the Intermediate Value Theorem follows that there exists an $\alpha \in (\alpha_1, \alpha_2)$ s.t. $|\langle \mathbf{h}(\alpha), \mathbf{e} \rangle| = 0$. Thus by definition $s(\mathbf{h}(\alpha)) = |\langle \mathbf{h}(\alpha), \mathbf{e} \rangle| / \|\mathbf{h}(\alpha)\|$, we see that $\min_\alpha s(\mathbf{h}(\alpha)) = 0$.

On the other hand, since σ_a is a non-trivial linear map for $x \gg 0$, we may assume $\sigma_a(x) = cx$ for $x > x_0$ where $c \neq 0$ is some non-zero constant and $x_0 > 0$ is some positive constant. Then we can choose an $\alpha_0 > \langle \mathbf{z}, \mathbf{e} \rangle$ s.t. for any $\alpha \geq \alpha_0$, all of the entries of $\mathbf{z}(\alpha)$ are greater than x_0 . Then whenever $\alpha \geq \alpha_0$, we have $\mathbf{h}(\alpha) = \sigma_a(\mathbf{z}(\alpha)) = c\mathbf{z}(\alpha)$. This implies

$$s(\mathbf{h}(\alpha)) = \frac{|\langle \mathbf{h}(\alpha), \mathbf{e} \rangle|}{\|\mathbf{h}(\alpha)\|} = \frac{|\langle c\mathbf{z}(\alpha), \mathbf{e} \rangle|}{\|c\mathbf{z}(\alpha)\|} = \frac{|\langle \mathbf{z}(\alpha), \mathbf{e} \rangle|}{\|\mathbf{z}(\alpha)\|} = s(\mathbf{z}(\alpha)).$$

Thus $\sup_\alpha s(\mathbf{h}(\alpha)) = 1$ follows from the Proof of Proposition B.1 where we see that $\sup_\alpha s(\mathbf{z}(\alpha)) = 1$ since $s(\mathbf{z}(\alpha))$ gets closer to 1 as α increases. \square

Remark B.2. Indeed, it holds for any continuous function $f : \mathbb{R} \rightarrow \mathbb{R}$ satisfying the following

1. $f(x) > 0$ for $x \gg 0$, $f(x) < 0$ for $x \ll 0$ or $f(x) < 0$ for $x \gg 0$, $f(x) > 0$ for $x \ll 0$,
2. f is a non-trivial linear map for $x \gg 0$ or $x \ll 0$.

One can check that the proof above only depends on these two properties. It is worth mentioning that most activation functions, e.g., leaky LU, SiLU, tanh, satisfy condition 1.

Proof of Corollary 4.5. For any α , we notice that $\|\mathbf{z}\|_{\mathcal{M}^\perp} = \|\mathbf{z}_{\mathcal{M}^\perp}\|_F = \|\mathbf{z}(\alpha)\|_{\mathcal{M}^\perp}$ since α only changes the component of \mathbf{z} in the eigenspace \mathcal{M} . Also, Propositions 3.2 and 3.3 show that $\|\mathbf{z}(\alpha)\|_{\mathcal{M}^\perp} \geq \|\mathbf{h}(\alpha)\|_{\mathcal{M}^\perp}$ whenever $\mathbf{h}(\alpha) = \sigma(\mathbf{z}(\alpha))$ or $\sigma_a(\mathbf{z}(\alpha))$. Therefore, $\|\mathbf{z}\|_{\mathcal{M}^\perp} \geq \|\mathbf{h}(\alpha)\|_{\mathcal{M}^\perp}$ holds for any α . Since $\mathbf{z}_{\mathcal{M}^\perp} \neq 0$, $s(\mathbf{z})$ must lie in $[0, 1)$. \square

C Additional Experimental Details

This section provides additional experimental details and results for Section 6. All tasks were run on Nvidia RTX 3090, GV100, and Tesla T4 GPUs. Computational performance metrics, including timing, were measured using Tesla T4 GPUs on Colab.

	# Nodes	# Edges	# Features	# Classes	Splits (Train/Val/Test)
Cornell	183	295	1,703	5	48/32/20%
Texas	181	309	1,703	5	48/32/20%
Wisconsin	251	499	1,703	5	48/32/20%
Chameleon	2,277	36,101	2,325	5	48/32/20%
Squirrel	5,201	217,073	2,089	5	48/32/20%
Citeseer	3,727	4,732	3,703	6	120/500/1000
Cora	2,708	5,429	1,433	7	140/500/1000
PubMed	19,717	44,338	500	3	60/500/1000
Coauthor-Physics	34,493	247,962	8415	5	100/150/34,243
Ogbn-arxiv	169,343	1,166,243	128	40	90,941/29,799/48,603

Table 4: Graph statistics.

C.1 Dataset details

Table 4 presents additional graph statistics.

Citation Datasets: We use Cora, Citeseer, PubMed, Coauthor-Physics, and Ogbn-arxiv, where each dataset is a graph with nodes representing academic publications, features as bag-of-words vectors, labels indicating publication types, and edges denoting citations.

Web Knowledge-Base Datasets: We use the Cornell, Texas, and Wisconsin datasets, each represented as a graph where nodes are CS department webpages, features are bag-of-words vectors, edges denote hyperlinks, and labels indicate webpage types.

Wikipedia Network Datasets: We use the Chameleon and Squirrel datasets, where each graph represents CS department webpages as nodes, with bag-of-words features, hyperlink edges, and labels indicating webpage types.

C.2 Model size and computational time for citation datasets

Table 5 compares the model size and computational time for experiments on citation datasets in Section 6.2.

	# Parameters	Training Time (s)	Inference Time (ms)
Cora			
GCN	100,423	8.4	1.6
GCNII	110,535	10.0	2.1
GCNII	708,743	57.6	12.3
GCNII-SCT	1,237,127	110.3	29.6
EGNN	712,839	65.6	14.4
EGNN-SCT	316,551	24.8	4.5
Citeseer			
GCN	245,638	8.3	1.5
GCN-SCT	301,830	15.5	4.0
GCNII	999,174	57.6	12.3
GCNII-SCT	1,001,222	65.9	15.7
EGNN	739,078	39.6	7.2
EGNN-SCT	540,934	24.0	5.8
PubMed			
GCN	40,451	9.0	1.8
GCN-SCT	40,707	11.1	2.2
GCNII	326,659	98.2	12.8
GCNII-SCT	590,851	71.7	17.4
EGNN	592,899	93.7	2.5
EGNN-SCT	130,563	16.0	3.1
Coauthor-Physics			
GCN	547,141	35.2	8.0
GCN-SCT	547,397	33.9	8.3
GCNII	555,333	49.1	10.3
GCNII-SCT	555,461	67.0	9.5
EGNN	672,069	176.4	47.9
EGNN-SCT	572,229	51.7	14.8
Ogbn-arxiv			
GCN	27,240	50.4	21.1
GCN-SCT	28,392	62.6	24.4
GCNII	76,392	205.4	94.8
GCNII-SCT	80,616	253.0	108.9
EGNN	77,416	206.8	98.0
EGNN-SCT	81,640	254.0	112.3

Table 5: Number of model parameters for varying numbers of layers using the optimal model hyperparameters. The SCT is added at each layer, and the size of the additional parameters scales with the number of eigenvectors with an eigenvalue of one for matrix \mathbf{G} in equation 2.

Parameter	Values
Learning Rate	{ $1e-4, 1e-3, 1e-2$ }
Weight Decay (FC)	{ $0, 1e-4, 5e-4, 1e-3, 5e-3, 1e-2$ }
Weight Decay (Conv)	{ $0, 1e-4, 5e-4, 1e-3, 5e-3, 1e-2$ }
Dropout	{ $0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9$ }
Hidden Channels	{ $16, 32, 64, 128$ }
GCNII- α	{ $0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9$ }
GCNII- θ	{ $0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9$ }
EGNN- c_{\max}	{ $0.5, 1.0, 1.5, 2.0$ }
EGNN- α	{ $0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9$ }
EGNN- θ	{ $0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9$ }

Table 6: Hyperparameter grid search for Table 1.

C.3 Additional Section 6.2 details for citation datasets

Table 6 lists the hyperparameters used in the grid search for generating the results in Table 1. Table 8 reports classification accuracy across model depths using ReLU and leaky ReLU activations.

Layers	2	4	16	32
Cora				
EGNN/EGNN-SCT	83.2/ 83.4	84.2/ 84.3	85.4/ 85.5	85.3/ 85.5
Citeseer				
EGNN/EGNN-SCT	72.0/ 72.1	71.9/ 72.3	72.4/ 72.6	72.3/ 72.8
PubMed				
EGNN/EGNN-SCT	79.2/ 79.4	79.5/ 79.8	80.1/80.1	80.0/ 80.2
Coauthor-Physics				
EGNN/EGNN-SCT	92.6/ 92.8	92.9/ 93.0	93.1/ 93.3	93.3/93.3
Ogbn-arxiv				
EGNN/EGNN-SCT	68.4/ 68.5	71.1/ 71.3	72.7/ 73.0	72.7/ 72.9

Table 7: Test accuracy for EGNN and EGNN-SCT using the SReLU activation function of varying depth on citation networks with the split discussed in Section 6.2. (Unit:%)

C.3.1 Vanishing gradients

Figure 4 illustrates the vanishing gradient problem in training deep GCNs—with and without SCT—compared to GCNII and EGNN. It plots $\|\partial \mathbf{H}^{\text{out}} / \partial \mathbf{H}^l\|$ for layers $l \in [0, 32]$ over 100 training epochs. Subfigures (a) and (b) show that GCN and GCN-SCT suffer from vanishing gradients. In contrast, (c) and (e) demonstrate that GCNII and EGNN avoid this issue by connecting \mathbf{H}^0 to every layer, maintaining nonzero gradients in early layers. Interestingly, adding SCT to GCNII and EGNN (shown in (d) and (f)) amplifies intermediate gradients as training progresses.

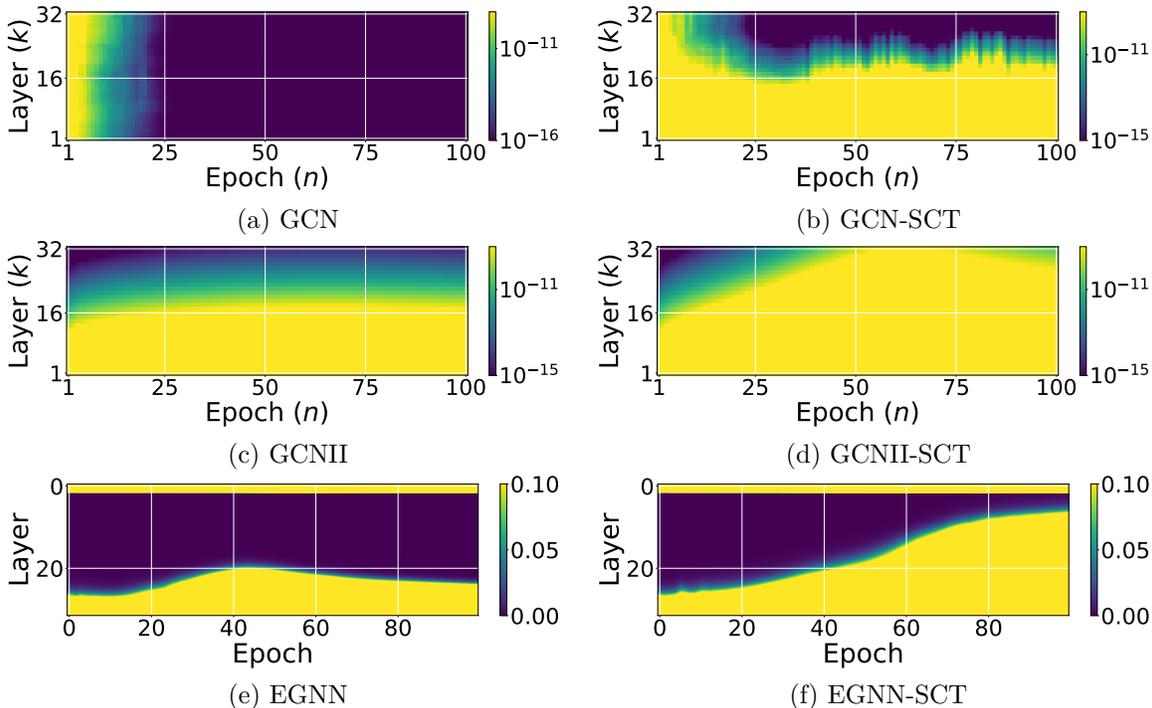


Figure 4: We visualize training gradients $\|\partial \mathbf{H}^{\text{out}} / \partial \mathbf{H}^l\|$ for layers $l \in [0, 32]$ over 100 epochs on the Citeseer dataset. All models use 32 layers with 16 hidden dimensions. Subfigure (a) shows that GCN suffers from vanishing gradients, while (c) GCNII and (e) EGNN maintain nonzero gradients due to skip connections to \mathbf{H}^0 . Although SCT does not resolve vanishing gradients in (b) GCN-SCT, it increases intermediate gradient norms in (d) GCNII-SCT and (f) EGNN-SCT as training progresses.

Cora								
	ReLU				leaky ReLU			
Layers	2	4	16	32	2	4	16	32
GCN-SCT	81.2	80.3	71.4	67.2	82.9	82.8	68.0	65.5
GCNII-SCT	83.5	83.8	82.7	83.3	83.8	84.8	84.8	85.5
EGNN-SCT	84.1	83.8	82.3	80.8	83.7	84.5	83.3	82.0
Citeseer								
	ReLU				leaky ReLU			
Layers	2	4	16	32	2	4	16	32
GCN-SCT	69.0	67.3	51.5	50.3	69.9	67.7	55.4	51.0
GCNII-SCT	72.8	72.8	72.8	73.3	72.8	72.9	73.8	72.7
EGNN-SCT	72.5	72.0	70.2	71.8	73.1	71.7	72.6	72.9
PubMed								
	ReLU				leaky ReLU			
Layers	2	4	16	32	2	4	16	32
GCN-SCT	79.4	78.2	75.9	77.0	79.8	78.4	76.1	76.9
GCNII-SCT	79.7	80.1	80.7	80.7	79.6	80.0	80.3	80.7
EGNN-SCT	79.7	80.1	80.0	80.4	79.8	80.4	80.3	80.2
Coauthor-Physics								
	ReLU				leaky ReLU			
Layers	2	4	16	32	2	4	16	32
GCN-SCT	91.8 ± 1.6	91.6 ± 3.0	44.5 ± 13.0	42.6 ± 17.0	92.6 ± 1.6	92.5 ± 5.9	50.9 ± 15.0	43.6 ± 16.0
GCNII-SCT	94.4 ± 0.4	93.5 ± 1.2	93.7 ± 0.7	93.8 ± 0.6	94.0 ± 0.4	94.2 ± 0.3	93.3 ± 0.7	94.1 ± 0.3
EGNN-SCT	93.6 ± 0.7	94.1 ± 0.4	93.4 ± 0.8	93.8 ± 1.3	93.9 ± 0.7	94.0 ± 0.7	94.0 ± 0.7	93.3 ± 0.9
Ogbn-arxiv								
	ReLU				leaky ReLU			
Layers	2	4	16	32	2	4	16	32
GCN-SCT	71.7 ± 0.3	72.6 ± 0.3	71.4 ± 0.2	71.9 ± 0.3	72.1 ± 0.3	72.7 ± 0.3	72.3 ± 0.2	72.3 ± 0.3
GCNII-SCT	71.4 ± 0.3	72.1 ± 0.3	72.2 ± 0.2	71.8 ± 0.2	72.0 ± 0.3	72.2 ± 0.2	72.4 ± 0.3	72.1 ± 0.3
EGNN-SCT	68.5 ± 0.6	71.0 ± 0.5	72.8 ± 0.5	72.1 ± 0.6	67.7 ± 0.5	71.3 ± 0.5	72.3 ± 0.5	72.3 ± 0.5

Table 8: Test accuracy for models of varying depth using ReLU or leaky ReLU activation on citation network datasets, based on the split described in Section 6.2.

C.4 Additional Section 6.2 details for other datasets

Table 9 reports mean test accuracy and standard deviation over 10 folds for SCT-based models on the WebKB and WikipediaNetwork datasets. Table 10 shows average per-epoch time for models with 8 layers. These results indicate that integrating SCT incurs minimal computational overhead.

	Cornell	Texas	Wisconsin	Chameleon	Squirrel
GCN-SCT	55.95 ± 8.5	62.16 ± 5.7	54.71 ± 4.4	38.44 ± 4.3	35.31 ± 1.9
GCNII-SCT	75.41 ± 2.2	83.34 ± 4.5	86.08 ± 3.8	64.52 ± 2.2	47.51 ± 1.4

Table 9: Mean ± standard deviation test accuracy from 10-fold cross-validation on five heterophilic datasets with fixed 48/32/20

	Cornell	Texas	Wisconsin	Chameleon	Squirrel
GCN Kipf & Welling (2017)	0.011	0.013	0.012	0.011	0.022
GCNII Chen et al. (2020b)	0.017	0.018	0.017	0.013	0.022
GCN-SCT	0.015	0.017	0.015	0.011	0.023
GCNII-SCT	0.017	0.018	0.017	0.020	0.025

Table 10: Average per-epoch computation time on five heterophilic datasets with fixed 48/32/20% splits. All models use 8 layers with 16 hidden channels. (Unit: second)