GRADIENT FLOW PROVABLY LEARNS ROBUST CLASSI FIERS FOR DATA FROM ORTHONORMAL CLUSTERS

Anonymous authors

Paper under double-blind review

ABSTRACT

Deep learning-based classifiers are known to be vulnerable to adversarial attacks. Existing methods for defending against such attacks require adding a defense mechanism or modifying the learning procedure (e.g., by adding adversarial examples). This paper shows that for certain data distribution one can learn a provably robust classifier using standard learning methods and without adding a defense mechanism. More specifically, this paper addresses the problem of finding a robust classifier for a binary classification problem in which the data comes from a mixture of Gaussian clusters with orthonormal cluster centers. First, we characterize the largest ℓ_2 -attack any classifier can defend against while maintaining high accuracy, and show the existence of optimal robust classifiers achieving this maximum ℓ_2 -robustness. Next, we show that given data sampled from the orthonormal cluster model, gradient flow on a two-layer network with a polynomial ReLU activation and without adversarial examples provably finds an optimal robust classifier.

1 INTRODUCTION

025 026

004

010 011

012

013

014

015

016

017

018

019

021

023 024

The vulnerability of deep neural networks to *adversarial attacks* (Szegedy et al., 2014), which are typically human-imperceptible perturbations to the input data, has led to numerous efforts in building defenses against these attacks (Shafahi et al., 2019; Papernot et al., 2016; Wong et al., 2019; Guo et al., 2018; Cohen et al., 2019; Levine & Feizi, 2020; Yang et al., 2020; Sulam et al., 2020; Kinfu & Vidal, 2022). These defenses have been counteracted by new adaptive attacks (Athalye et al., 2018; Carlini et al., 2019; Croce & Hein, 2020), leading to new defenses and so on. Even in the era of Large Language Models, adversarial attacks exist (Chao et al., 2023; Shah et al., 2023), leading to undesired or harmful model outputs, and the competition between adversaries and defenders continues (Robey et al., 2023; Ji et al., 2024). While such a competition allows us to design more robust networks, it will not end unless many fundamental questions about adversarial robustness are answered.

One question is what is the maximum adversarial perturbation a neural network can tolerate? Many 037 works on certified robustness (Cohen et al., 2019; Fazlyab et al., 2020; Zhang et al., 2018) aim to 038 find a certified radius such that a neural network can provably maintain a high prediction accuracy for adversarial attacks within that radius. However, their reported certified radii are often too small compared to what can be achieved by practical defenses (Tramèr et al., 2018; Guo et al., 2018; Gowal 040 et al., 2020; Wu et al., 2020). Yet, practical defenses come at the cost of computing adversarial 041 examples, or sophisticated model designs, mostly without theoretical guarantees, except for the case 042 of linear classifiers (Zou et al., 2021). This also motivates an intriguing question: Is it possible to 043 (provably) find a robust network by standard training methods, without adversarial examples? 044

We argue that these questions can be answered by exploiting properties of the data distribution, which most aforementioned works fail to do. Indeed, recent works show that the existence of a robust classifier is closely related to data geometry. For instance, Pal et al. (2023; 2024) show that if the *data is localized*, i.e., if the distribution of the data given the class concentrates in a set of small volume, then a robust classifier is guaranteed to exist. Moreover, they show that a 2r separation (w.r.t. to some distance metric) between the sets that contain each class-conditioned probability mass is sufficient for the existence of a robust classifier against attacks of radius r in the same distance metric.

This paper shows that such a relationship between data geometry and adversarial robustness has deeper implications: For certain data distributions, one can characterize the maximum robustness any classifier can achieve, based on how class-conditional probability masses are separated. Moreover, 054

056

058

060

061

062 063

064

065

066

067

068

069

070 071

072

074

075

076

077

078

079

080

081

082

084

090

091

092

094

one can make suitable architectural designs that exploit the data geometry, such that a nearly optimal robust classifier is provably learned by standard training methods, such as gradient descent. Specifically, we consider a balanced mixture of K-Gaussian clusters in \mathbb{R}^D , split into two classes:

$$\underbrace{\mathcal{N}\left(\boldsymbol{\mu}_{1}, \alpha^{2}\boldsymbol{I}/D\right), \cdots, \mathcal{N}\left(\boldsymbol{\mu}_{K_{1}}, \alpha^{2}\boldsymbol{I}/D\right)}_{\text{positive (+1) class}}, \underbrace{\mathcal{N}\left(\boldsymbol{\mu}_{K_{1}+1}, \alpha^{2}\boldsymbol{I}/D\right), \cdots, \mathcal{N}\left(\boldsymbol{\mu}_{K}, \alpha^{2}\boldsymbol{I}/D\right)}_{\text{negative (-1) class}}, \quad (1)$$

where the *cluster centers* $\mu_1, \dots, \mu_K \in \mathbb{R}^D$ are othonormal, α^2 denotes the *intra-cluster variance*, and the ambient dimension D is sufficiently large. We explain our contributions as follows:

Maximum ℓ_2 -robustness This mixture of Gaussian distribution satisfies data localization and separation properties similar to those studied in Pal et al. (2023). As illustrated in Figure 1 for the case of two clusters (one from the positive class and one from the negative class), the class-conditioned probability masses concentrate around two (D-1)-dimensional affine subspaces separated by a Euclidean distance of almost $\sqrt{2}$. Based on such observation, our first set of results are:

Theorem (Theorem 1 & 2, informal). No classifier can defend against an adversarial attack of ℓ_2 radius $\frac{\sqrt{2}}{2}$. However, one can construct a nearly optimal robust classifier that can defend against attacks of radius arbitrarily close to $\frac{\sqrt{2}}{2}$ when D is sufficiently large.

Our results show that data localization and separation are important properties in understanding the maximum achievable robustness for a classifier. Moreover, we will show that the classifier we construct is the Bayes optimal classifier w.r.t. the 0-1 loss, which operates as a nearest-cluster rule: classifiers that exploit the multi-cluster data structure are naturally and optimally robust.



Figure 1: Illustration of two clusters in high-dimensions, each concentrated on a (D-1)-dimensional affine subspace such that the subspaces are separated by a Euclidean distance of $\sqrt{2}$.



Figure 2: Given sampled data from (1) with 12 positive clusters and 8 negative clusters (D = 2000), gradient descent (SGD, small initialization) on (bias-free, width-200) two-layer ReLU network (ReLU) fails to find a robust classifier. This issue persists after 1) increasing depth to 4 (MLP); 2) (blindly) switching to another activation (Tanh); or 3) using a linear classifier (LogReg). However, by choosing a suitable activation (pReLU, p = 3), GD can find a nearly optimal robust classifier.

096 **Learning optimal robust networks** So far everything seems to be intuitive and straightforward given the fairly simple distributional assumption. However, issues arise when one does not know 098 the data distribution a priori and seeks a classifier by training a neural network on sampled data via 099 gradient descent. As Figure 2 suggests, a trained multi-layer ReLU network fails to find a classifier with the same level of robustness as the Bayes classifier (which indeed can defend against attack of 100 101 radius $\sim \frac{\sqrt{2}}{2}$, as our results suggest). This matter is first discussed by Frei et al. (2023), where they show that any two-layer ReLU network trained by gradient descent under data samples from (1) is 102 non-robust against adversarial attacks of ℓ_2 -radius $\Theta(\frac{1}{\sqrt{K}})$, where K is the total number of clusters. 103 Later, Min & Vidal (2024) show that this issue is caused by the fact that a ReLU network fails to 104 105 learn, internally with its weight parameters, the multi-cluster structure of the data distribution, despite that the sampled data points are revealing such a structure. Therefore, while the structural property of 106 the data distribution allows one to construct an optimal robust classifier, gradient descent algorithms 107 on neural networks may struggle to learn these key properties, leading to non-robust classifiers.

To address this issue, Min & Vidal (2024) propose to change the activation. More specifically, replacing the ReLU activation with a polynomial ReLU activation (pReLU) with polynomial degree p as a hyperparameter. They empirically show that when p = 3, the pReLU network can internally learn the data structure, leading to a more robust classifier, and they conjecture that this improvement in robustness happens when $p \ge 3$. However, a rigorous analysis of convergence is not provided. Our second set of results is to develop a full convergence analysis for gradient flow, a continuous time limit of gradient descent by taking the stepsize to zero, on a two-layer pReLU network and show that:

Theorem (Theorem 3 & Corollary 1, informal). When p > 2 and the intra-cluster variance α^2 is sufficiently small, gradient flow on pReLU networks converges to a nearly optimal robust classifier.

117

Our analysis is based on prior works on gradient descent/flow with small initialization on two-layer ReLU networks Maennel et al. (2018); Phuong & Lampert (2021); Boursier et al. (2022); Kumar & Haupt (2024); Chistikov et al. (2023); Wang & Ma (2023); Min et al. (2024) and extends to pReLU networks. We show how the implicit bias of the gradient flow dynamics critically depends on a careful choice of activation function, allowing the network to learn accurately the underlying data structure, which, as we have discussed, is essential for finding a robust classifier.

123

Notation We denote the inner product between vectors \boldsymbol{x} and \boldsymbol{y} by $\langle \boldsymbol{x}, \boldsymbol{y} \rangle = \boldsymbol{x}^{\top} \boldsymbol{y}$, and the cosine of the angle between them as $\cos(\boldsymbol{x}, \boldsymbol{y}) = \langle \frac{\boldsymbol{x}}{\|\boldsymbol{x}\|}, \frac{\boldsymbol{y}}{\|\boldsymbol{y}\|} \rangle$. For an $n \times m$ matrix \boldsymbol{A} , we let $\|\boldsymbol{A}\|$ and $\|\boldsymbol{A}\|_F$ denote the spectral and Frobenius norm of \boldsymbol{A} , respectively. We also define $\mathbb{1}_A$ as the indicator for a statement A: $\mathbb{1}_A = 1$ if A is true and $\mathbb{1}_A = 0$ otherwise. We also let $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}^2)$ denote the normal distribution with mean $\boldsymbol{\mu}$ and covariance matrix $\boldsymbol{\Sigma}^2$, and Unif(S) denote the uniform distribution over a set S. Lastly, we let [N] denote the integer set $\{1, \dots, N\}$.

130 131

132

137 138

2 OPTIMAL ROBUST CLASSIFIERS FOR ORTHONORMAL CLUSTERS

Orthonormal cluster model We study a balanced mixture of K Gaussian clusters, and K_1 of them belong to the positive (+1) class and $K_2 := K - K_1$ of them the negative (-1) class. Formally, consider a tuple of random variables (X, Y, Z) on $\mathbb{R}^D \times \{+1, -1\} \times [K]$ representing *observed data*, *observed class label*, and *latent cluster membership*, respectively, defined as follow:

$$Z \sim \text{Unif}(\{1, \cdots, K\}), \ X|Z \sim \mathcal{N}\left(\boldsymbol{\mu}_{Z}, \alpha^{2} \boldsymbol{I}/D\right), \ Y|Z = \mathbb{1}_{Z \leq K_{1}} - \mathbb{1}_{Z > K_{1}},$$
(2)

where the μ_1, \dots, μ_K , called *cluster centers*, are a set of orthonormal vectors in \mathbb{R}^D , i.e. $\langle \mu_k, \mu_l \rangle = \mathbb{1}_{l=k}$. We denote the marginal distribution of (X, Y)-pair as $\mathcal{D}_{X,Y}$.

142 ℓ_2 -robust classifier for $\mathcal{D}_{X,Y}$ Our interest is to find a classifier that not only accurately predicts the label y given an observed data x, but do so in a way that is robust to some adversarial attacks on 143 observation x. Specifically, we search for a classifier $f : \mathbb{R}^D \to \mathbb{R}$ such that with high probability 144 $\min_{\|\boldsymbol{d}\|=1} f(\boldsymbol{x} + r\boldsymbol{d})y > 0$ for some $r \ge 0$ given a new sample (\boldsymbol{x}, y) from $\mathcal{D}_{X,Y}$. When r = 0, 145 f(x)y > 0 suggest that sign (f(x)) correctly predicts the label y; When r > 0, min_{||d||=1} f(x + 1)146 rd)y > 0 suggest that sign (f(x + rd)) still makes correct prediction on y even though observation 147 x has been corrupted by some adversarial attack rd, thus robust to adversarial attacks of ℓ_2 -norm 148 radius r. Ideally, we want a classifier that is robust to attack of radius r, with as large r as possible. 149

150 151 152 153 154 155 156 Maximum achievable ℓ_2 -robustness Inevitably, any classifier fails to be robust if the adversary 154 has too much power, i.e., the attack radius r exceeds some value. Indeed, for the data distribution 153 $\mathcal{D}_{X,Y}$ of our interest, no classifier can defend against attacks of radius $\frac{\sqrt{2}}{2}$, as formally shown below: 154 Theorem 1. Let $f : \mathbb{R}^D \to \mathbb{R}$ be any Lebesgue measurable function such that the random variable 155 $\min_{\|d\|\leq 1} \left[f\left(x + \frac{\sqrt{2}}{2}d\right)y \right]$ is also measurable. Given a sample $(x, y) \sim \mathcal{D}_{X,Y}$, we have

$$\mathbb{P}_{(\boldsymbol{x},\boldsymbol{y})\sim\mathcal{D}_{X,Y}}\left(\min_{\|\boldsymbol{d}\|\leq 1}\left[f\left(\boldsymbol{x}+\frac{\sqrt{2}}{2}\boldsymbol{d}\right)\boldsymbol{y}\right]\leq 0\right)\geq\frac{\min\{K_1,K_2\}}{K}.$$
(3)

158 159

157

We refer the readers to Appendix B.1 for the proof. We explain Theorem 1 from a geometric perspective (we have discussed some in the introduction): Consider the case of two clusters $\mathcal{N}(\boldsymbol{\mu}_1, \frac{\alpha^2}{D}\boldsymbol{I})$ and $\mathcal{N}(\mu_2, \frac{\alpha^2}{D}I)$ of different classes. As shown in Figure 1, when ambient dimension D is large, we expect that each cluster concentrates around a D-1 affine subspace that is orthogonal to the vector $\mu_1 - \mu_2$. Most importantly, the distance between these two affine subspaces is $\sqrt{2}$, suggesting that given any decision boundary that separates two affine subspaces, an adversary can perturb a substantial portion of the probability mass of these clusters to cross the boundary with an attack radius $\frac{\sqrt{2}}{2}$. The same argument holds for the K-clusters cases, where every two clusters are separated by a Euclidean distance $\sqrt{2}$. We also note that extending Theorem 1 to attacks in another metric amounts to measuring this separation in that metric. Our second result shows the Bayes optimal classifier w.r.t. 0-1 loss is also nearly optimally robust:

Theorem 2. The Bayes optimal classifier for label Y given observation \mathbf{x} w.r.t. 0-1 loss is sign $(f^*(\mathbf{x}))$, where $f^*(\mathbf{x}) = \sum_{k=1}^{K_1} \exp\left(\frac{D\langle \mathbf{x}, \boldsymbol{\mu}_k \rangle}{\alpha^2}\right) - \sum_{k=K_1+1}^{K} \exp\left(\frac{D\langle \mathbf{x}, \boldsymbol{\mu}_k \rangle}{\alpha^2}\right)$. Moreover, given a sample $(\mathbf{x}, y) \sim \mathcal{D}_{X,Y}$, we have, for any $\frac{2\sqrt{2}\alpha^2 \log K}{D} \le \nu \le \sqrt{2}$,

$$\mathbb{P}_{(\boldsymbol{x},\boldsymbol{y})\sim\mathcal{D}_{\boldsymbol{X},\boldsymbol{Y}}}\left(\min_{\|\boldsymbol{d}\|\leq 1}\left[f^*\left(\boldsymbol{x}+\frac{\sqrt{2}-\nu}{2}\boldsymbol{d}\right)\boldsymbol{y}\right]>0\right)\geq 1-2K\exp\left(-\frac{D\nu^2}{64\alpha^2}\right).$$
 (4)

We refer the readers to Appendix B.2 for the proof. If we pick $\nu = \Theta(\left(\frac{\alpha^2}{D}\right)^{\frac{1}{4}})$ in Theorem 2, then the result shows that f^* is robust against attacks of radius $\frac{\sqrt{2}}{2} - \Theta(\left(\frac{\alpha^2}{D}\right)^{\frac{1}{4}})$ with probability at least $1 - \mathcal{O}(K \exp\left(-\left(\frac{D}{\alpha^2}\right)^{\frac{1}{2}}\right))$ over new sample from $\mathcal{D}_{X,Y}$. Therefore, f^* is nearly optimal robust when $\frac{\alpha^2}{D} = o(1)$, i.e. the ambient dimension is large or the intra-class variance is small.

Interpreting f^* as a nearest-cluster rule We explain why this Bayes classifier is of interest. We have the following derivation:

$$\operatorname{sign}\left(f^{*}(\boldsymbol{x})\right) = \operatorname{sign}\left(\sum_{k=1}^{K_{1}} \exp\left(\frac{D\left\langle\boldsymbol{x},\boldsymbol{\mu}_{k}\right\rangle}{\alpha^{2}}\right) - \sum_{k=K_{1}+1}^{K} \exp\left(\frac{D\left\langle\boldsymbol{x},\boldsymbol{\mu}_{k}\right\rangle}{\alpha^{2}}\right)\right)$$

$$= \operatorname{sign}\left(\frac{\alpha^2}{D}\log\left(\sum_{k=1}^{K_1} \exp\left(\frac{D\langle \boldsymbol{x}, \boldsymbol{\mu}_k \rangle}{\alpha^2}\right)\right) - \frac{\alpha^2}{D}\log\left(\sum_{k=K_1+1}^{K} \exp\left(\frac{D\langle \boldsymbol{x}, \boldsymbol{\mu}_k \rangle}{\alpha^2}\right)\right)\right)$$
$$= \operatorname{sign}\left(-\max_{k=1}^{K} \langle \boldsymbol{x}, \boldsymbol{\mu}_k \rangle - \max_{k=1}^{K} \langle \boldsymbol{x}, \boldsymbol{\mu}_k \rangle + \mathcal{O}\left(\log K\frac{\alpha^2}{\alpha^2}\right)\right)$$

$$= \operatorname{sign}\left(\max_{1 \le k \le K_1} \langle \boldsymbol{x}, \mu_k \rangle - \max_{K_1 + 1 \le k \le K} \langle \boldsymbol{x}, \mu_k \rangle + \mathcal{O}\left(\log K \frac{\alpha}{D}\right)\right)$$

where the second inequality is due to the fact that $\frac{\alpha}{D} \log(\cdot)$ function is a non-decreasing function, and the third inequality is because LogSumExp $(\{z_1, \dots, z_K\})$ function with a temperature $\frac{D}{\alpha^2}$ uniformly approximate max_k z_k with an error $\mathcal{O}(\log K \frac{\alpha^2}{D})$. When the error is small, the Bayes classifier $f^*(x)$ finds the closest cluster center to x and outputs the label to that cluster, which is a nearest-cluster rule. Therefore, by exploiting the multi-cluster structure of $\mathcal{D}_{X,Y}$, f^* achieves the maximum ℓ_2 -robustness.

So far we have shown that a nearly optimal robust classifier for $\mathcal{D}_{X,Y}$ can be easily constructed as a nearest-cluster rule. However, as we discussed in the introduction, gradient descent algorithms with sampled data often fail to find a classifier with the same level of robustness. Next, we address the problem of finding a nearly optimal robust classifier by gradient flow dynamics.

3 Optimal Robust Classifiers Obtained via Gradient Flow

In this section, we aim to find a nearly optimal ℓ_2 -robust classifier for $\mathcal{D}_{X,Y}$ by vanilla gradient descent without adversarial training. We start by stating the problem of training two-layer networks with gradient flow (gradient descent with infinitesimal step size). Then we show that with a pReLU activation, gradient flow provably finds a classifier that is nearly optimal ℓ_2 -robust.

3.1 PRELIMINARIES: GRADIENT FLOW ON TWO-LAYER NETWORKS

pReLU network We consider a two-layer pReLU network (Min & Vidal, 2024) defined as follow:

$$f^{(p)}(\boldsymbol{x};\boldsymbol{\theta}) = \sum_{i=1}^{h} v_j \frac{\sigma^{p}(\langle \boldsymbol{x}, \boldsymbol{w}_j \rangle)}{\|\boldsymbol{w}_j\|^{p-1}} \qquad (\boldsymbol{\theta} := \{\boldsymbol{w}_j, v_j\}_{j=1}^{h}),$$
(5)

216 $f^{(p)}$ can be viewed as a generalized version of the ReLU network. When p = 1, $f^{(1)}$ is exactly a 217 two-layer ReLU network. When p > 1, the output of the hidden activation is equal to the one of 218 the ReLU network multiplying $\cos^{p-1}(x, w_j)$ (Min & Vidal, 2024), which discourages large angle 219 separation between data x and *neuron* w_j .

220

229

230

231

232 233 234

239

240

241

251

252

 $\begin{array}{l} \begin{array}{l} \begin{array}{l} \begin{array}{l} 221\\ \ell_2\text{-loss function and balanced dataset} & \text{Given a dataset } \{x_i, y_i\}_{i=1}^n, \text{ one define the loss function as} \\ \begin{array}{l} \mathcal{L}(\theta; \{x_i, y_i\}_{i=1}^n) = \sum_{i=1}^n \ell(y_i, \hat{y}_i), \text{ where } \hat{y}_i = f^{(p)}(x_i; \theta) \text{ . For classification problem, the typical} \\ \begin{array}{l} \text{choice of } \ell \text{ can be exponential } \exp(-y\hat{y}), \text{ or logistic loss } \log(1 + \exp(-y\hat{y})). \end{array} \end{array} \right. \text{Most of our theoretical} \\ \begin{array}{l} \text{analysis works for these choices for } \ell. \end{array} \\ \begin{array}{l} \text{However, using classification losses poses additional challenges} \\ \text{in analyzing the late phase of the training (details explained in later sections). Therefore, our theorem \\ \text{considers a } \ell_2\text{-loss: } \ell(y, \hat{y}) = \frac{1}{2} \|y - \hat{y}\|^2, \text{ and the extension to classification losses is discussed in} \\ \text{Section 3.2.4.} \end{array}$

As for the dataset, since $\mathcal{D}_{X,Y}$ samples data with equal probability from each cluster, there are approximately equal number of samples from each cluster when we sample a large number of data. Therefore, instead of considering a dataset directly sampled from $\mathcal{D}_{X,Y}$, we consider the following balanced dataset $\hat{\mathcal{D}} = \{x_i, y_i\}_{i=1}^{KN}$, where

$$\boldsymbol{x}_{i} \sim \mathcal{N}\left(\boldsymbol{\mu}_{k}, \alpha^{2}\boldsymbol{I}/\boldsymbol{D}\right), y_{i} = \mathbb{1}_{k \leq K_{1}} - \mathbb{1}_{k > K_{1}}, \quad (k-1)N+1 \leq i \leq kN, \ 1 \leq k \leq K.$$
(6)

We call this dataset balanced because \hat{D} has exactly N samples from each cluster $\mathcal{N}(\mu_k, \alpha^2 I/D)$. This assumption allows us to omit the additive perturbations in our analysis introduced by unbalanced per-cluster sample size.

Gradient flow with small and balanced initialization Given the network parametrization θ and the loss function \mathcal{L} constructed from a balanced dataset $\hat{\mathcal{D}}$, we consider training the network by the following *gradient flow* (GF) dynamics¹:

$$\dot{\boldsymbol{\theta}} = -\nabla_{\boldsymbol{\theta}} \mathcal{L}\left(\boldsymbol{\theta}; \hat{\mathcal{D}}\right), \qquad \boldsymbol{\theta}(0) = \boldsymbol{\theta}_0, \qquad (7)$$

We assume the initialization $\theta(0)$ is ϵ -small and balanced, formally defined as the following.

Assumption 1 (ϵ -small and balanced initialization). The initialization $\boldsymbol{\theta}(0) = \{\boldsymbol{w}_j(0), v_j(0)\}_{j=1}^h$ satisfies the following: there exists an initialization shape $\{\boldsymbol{w}_{j0}, v_{j0}\}_{j=1}^h$ with $W_{\min} \leq \|\boldsymbol{w}_{j0}\| \leq W_{\max}, \forall j$, for some $W_{\min}, W_{\max} > 0$ and an initialization scale $\epsilon > 0$ such that

$$\boldsymbol{w}_{j}(0) = \epsilon \boldsymbol{w}_{j0}, \ v_{j}(0) = \epsilon v_{j0}, \ \|\boldsymbol{w}_{j0}\| = |\boldsymbol{v}_{j0}|, \forall j.$$
 (8)

253 Under a balanced initialization, we have $\|\boldsymbol{w}_i(0)\| = |v_i(0)|, \forall j$, and this balancedness holds through-254 out GF trajectory (See Appendix D.1): $||w_j(t)|| = |v_j(t)|, \forall j$. The balancedness between w_j and 255 v_i allows us to focus on the dynamics of w_i , which has been a common assumption in prior work 256 of this type (Maennel et al., 2018; Boursier et al., 2022; Chistikov et al., 2023; Min et al., 2024). Readers may view this assumption as made out of convenience, but it is essential for a tractable 257 analysis (also allowing an elegant interpretation of dynamics of w_i (Maennel et al., 2018; Boursier & 258 Flammarion, 2024)), and the theoretical results out of this assumption match the empirical results 259 when no balancedness is enforced (Min et al., 2024). 260

Given a balanced initialization, one can show that $sign(v_j(t)) = sign(v_j(0)), \forall j, \forall t \ge 0$ (Boursier et al., 2022). Roughly speaking, $sign(v_j(0))$ determines the dynamical behavior of neuron w_j under gradient flow: neurons with $sign(v_j(0)) = +1$ tend to align its direction with one of the positive cluster centers, μ_k , $k = 1, \dots, K_1$, and those with $sign(v_j(0)) = -1$ tend to align with one of the negative cluster centers. For this reason, we define the following neuron index sets: $\mathcal{N}_+ := \{j \in [h] : sign(v_j(0)) = +1\}$ and $\mathcal{N}_- := \{j \in [h] : sign(v_j(0)) = -1\}$.

¹Readers may find it more appropriate to study gradient flow as differential inclusion, instead of differential equation, since ReLU is non-differentiable at 0. However, our focus is on pReLU network with p > 2, which renders the network $f^{(p)}$ differentiable everywhere.

3.2 MAIN RESULTS: PRELU (p > 2) provably finds (NEAR)-optimal robust classifiers

pReLU classifier and the conjecture Min & Vidal (2024) study the adversarial robustness of the following pReLU classifier (that can be expressed by $f^{(p)}(x; \theta)$ with some choice of θ):

$$F^{(p)}(\boldsymbol{x}) = \sum_{k=1}^{K_1} \sigma^p(\langle \boldsymbol{x}, \boldsymbol{\mu}_k \rangle) - \sum_{k=K_1+1}^{K} \sigma^p(\langle \boldsymbol{x}, \boldsymbol{\mu}_k \rangle), \qquad (9)$$

and show that $F^{(p)}(x)$ is robust to adversarial attacks of ℓ_2 radius arbitrarily close to $\frac{\sqrt{2}}{2}$ when $\frac{D}{\alpha^2}$ is large (Now based on our Section 2, we know that $F^{(p)}(x)$ is nearly optimally robust). They conjecture that when p > 3 and the intra-cluster variance α^2 is small, the gradient flow on pReLU network $f^{(p)}(\cdot; \theta)$ with small initialization finds a classifier that is close to $F^{(p)}(x)$ up to a constant scaling factor. Then they argue that such proximity to $F^{(p)}(x)$ implies that the trained network has the same level of robustness as $F^{(p)}(x)$. Our main results fully prove this conjecture with p > 2.

Closeness to $F^{(p)}$ **implies robustness** We first show that given any classifier f(x) that is positively homogeneous of degree 1 w.r.t. x and is close to $F^{p}(x)$ in terms of some distance measure, it is nearly optimal robust when the intra-class variance is small (We refer to Appendix C for the proof.).

Proposition 1. Given a classifier f that satisfies $f(\gamma x) = \gamma f(x), \forall x \in \mathbb{R}^D, \forall \gamma > 0$ and $\operatorname{dist}(f, F^{(p)}) = \operatorname{inf}_{c>0} \sup_{\boldsymbol{x} \in \mathbb{S}^{D-1}} |cf(\boldsymbol{x}) - F^{(p)}(\boldsymbol{x})| \le \nu \text{ for some } p > 2 \text{ and } 0 < \nu \le \left(\frac{\sqrt{2}}{8}\right)^p.$ Then for a sample $(\mathbf{x}, y) \sim \mathcal{D}_{X,Y}$, we have

$$\mathbb{P}_{(\boldsymbol{x},\boldsymbol{y})\sim\mathcal{D}_{X,Y}}\left(\min_{\|\boldsymbol{d}\|\leq 1}\left[f\left(\boldsymbol{x}+\frac{\sqrt{2}-8\nu^{\frac{1}{p}}}{2}\boldsymbol{d}\right)\boldsymbol{y}\right]>0\right)\geq 1-2K\exp\left(-\frac{D\nu^{\frac{2}{p}}}{2K^{2}\alpha^{2}}\right)-4\exp\left(-\frac{3}{8\alpha^{2}}\right)$$
(10)

Given this result, it remains to show that gradient flow finds a network $f^{(p)}(\cdot; \theta)$ (which is positively homogeneous of degree 1) that is close to $F^{(p)}$ in the distance measure defined above. We will first discuss an additional assumption required on the initialization, then state our main result.

3.2.1 NON-DEGENERATE INITIALIZATION SHAPE

To properly define a non-degenerate initialization shape, we need to define a radial Voronoi tessella*tion* of $\mathbb{R}^{D-1}/\{0\}$ given a tuple of unit-norm vectors $\{\mu_k\}_{k\in\mathcal{K}}$.

Definition 1. Given a tuple of unit-norm vectors $\{\mu_k\}_{k \in \mathcal{K}}$, define the following $([\cdot]_+ := \max\{\cdot, 0\})$:

$$\mathcal{R}_k := \left\{ \boldsymbol{w} \in \mathbb{R}^{D-1} / \{0\} \mid [\cos(\boldsymbol{\mu}_k, \boldsymbol{w})]_+ > [\cos(\boldsymbol{\mu}_l, \boldsymbol{w})]_+, \forall l \neq k \right\}, k \in \mathcal{K},$$
 (Voronoi regions)
$$\mathcal{R}^\circ := \left\{ \boldsymbol{w} \in \mathbb{R}^{D-1} / \{0\} \mid [\cos(\boldsymbol{\mu}_k, \boldsymbol{w})]_+ = 0, \forall k \in \mathcal{K} \right\}.$$
(Void region)

From this definition, it is clear that $\{\mathcal{R}_k\}_{k\in\mathcal{K}}, \mathcal{R}^\circ$ are disjoint subsets of $\mathbb{R}^{D-1}/\{0\}$. We are ready to define a non-degenerate initialization shape, whose formal definition is stated below:

Definition 2 (Non-degenerate initialization shape). A set of initialization shape $\{w_{i0}\}_{i \in \mathcal{N}}$ is non*degenerate* w.r.t. a set of unit-norm vectors $\{\mu_k\}_{k \in \mathcal{K}}$ if it satisfies that

• (Neurons must be within one of the regions) $\forall j \in \mathcal{N}, w_{j0} \in (\bigcup_{k \in \mathcal{K}} \mathcal{R}_k) \bigcup \mathcal{R}^\circ$;

• (Non-void regions must contain at least one neuron) $\forall k \in \mathcal{K}, \exists j \in \mathcal{N}$ such that $w_{i0} \in \mathcal{R}_k$,

where $\{\mathcal{R}_k\}_{k\in\mathcal{K}}$ and \mathcal{R}° are the Voronoi regions and void region defined in Definition 1 w.r.t. $\{\mu_k\}_{k\in\mathcal{K}}$. Moreover, we let $d(w,S) = 1 - \sup_{s\in S, s\neq 0} \cos(w,s)$ and define non-degeneracy gap:

$$\Delta := \min\left\{\min_{\{\boldsymbol{w}_{j0}\in(\bigcup_{k\in\mathcal{K}}\mathcal{R}_k)\}} d\left(\boldsymbol{w}_{j0},\partial\left(\bigcup_{k\in\mathcal{K}}\mathcal{R}_k\right)\right),\min_{\{\boldsymbol{w}_{j0}\in\mathcal{R}^\circ\}} d\left(\boldsymbol{w}_{j0},\partial\mathcal{R}^\circ\right)\right\}.$$
 (11)

Whenever a vector w falls into one of the \mathcal{R}_k , it means that: 1) the angle between w and the corresponding μ_k is less than $\frac{\pi}{2}$; and 2) compared to all other μ_s , μ_k is the closest (in angle) to w. We hope that neurons initialized within some \mathcal{R}_k converge to the corresponding μ_k under GF, and those initialized within \mathcal{R}° stay in \mathcal{R}° (This is indeed the case, see Section 3.2.3).

330 The special case when a neuron is exactly initialized on the boundary of these Voronoi regions $\partial (\bigcup_{k \in \mathcal{K}} \mathcal{R}_k)$ cannot be analyzed since if a neuron 332 has equal angular distance to two μ vectors, there is no way to determine 333 which μ vector it converges to under GF with sampled data around these 334 μ vectors. Similarly, if a neuron is initialized at the boundary between some \mathcal{R}_k and \mathcal{R}° , then we can not determine whether it converges to μ_k , 335 or it falls into the interior of \mathcal{R}° and stays after that. Therefore we require 336 an initialization shape with a positive non-degeneracy gap. Moreover, 337 every \mathcal{R}_k must contain one neuron, ensuring the corresponding μ_k gets 338 learned. This leads to our assumption of non-degenerate initialization. 339



Figure 3: Illustration of a non-degenerate initialization shape $\{w_{10}, w_{20}\}$ w.r.t. two orthonormal vectors $\{\mu_1, \mu_2\}$.

Assumption 2. (Initialization has at least Δ non-degeneracy gap) $\exists \Delta > 0$ such that $\{w_{j0}\}_{j \in \mathcal{N}_{+}}$ is non-degenerate w.r.t. $\{\mu_k\}_{1 \leq k \leq K_1}$ with at least Δ non-degeneracy gap, and $\{w_{j0}\}_{j \in \mathcal{N}_{-}}$ is non-degenerate w.r.t. $\{\mu_k\}_{K_1 \leq k \leq K}$ with at least Δ non-degeneracy gap.

As one can see, this condition is stated per class: Positive (Negative) neurons must be initialized to be non-degenerate w.r.t. cluster centers from the positive (negative) class. We let $\{\mathcal{R}_k\}_{1 \le k \le K_1}$ and $\{\mathcal{R}_k\}_{K_1 \le k \le K}$ be the Voronoi regions defined by $\{\mu_k\}_{1 \le k \le K_1}$ and $\{\mu_k\}_{K_1 \le k \le K}$ respectively and define the neuron index sets $\mathcal{N}_k := \begin{cases} j \in \mathcal{N}_+ : w_{j0} \in \mathcal{R}_k, & 1 \le k \le K_1 \\ j \in \mathcal{N}_- : w_{j0} \in \mathcal{R}_k, & K_1 + 1 \le k \le K \end{cases}$ and $\mathcal{N}_c := [h] - \bigcup_{1 \le k \le K} \mathcal{N}_k$. As suggested in our previous discussion, we show that (See Section 3.2.3) under GF, all neurons in \mathcal{N}_k converge in angle to μ_k , which is an essential part of our theoretical results.

350 351

352

3.2.2 CONVERGENCE OF PRELU (p > 2) ON ORTHONORMAL CLUSTERS

Now we are ready to state our main theorem:

t

354 **Theorem 3** (pReLU converges to optimal robust classifier for orthonormal clusters). Let p > 2. Given 355 $0 \le \delta \le 1$ and a sufficiently small α_0^2 , consider data dimension $D \ge \tilde{\Omega}(\alpha_0^{-2})$ and per-cluster sample 356 size $\tilde{\Omega}(\alpha_0^{-2}) \leq N \leq \tilde{o}(\exp(\alpha_0^{-2}))$. With probability at least $1 - \delta$, the GF dynamics with a balanced 357 dataset $\hat{D} = \{x_i, y_i\}_{i=1}^{KN}$ sampled with intra-cluster variance $\alpha^2 \leq \alpha_0^2$, starting from some ϵ -small and balanced (Assumption 1) initialization $\theta(0)$ that satisfies Assumption 2 with a non-degeneracy 358 359 gap $\Delta = \Theta(1)$ and has a sufficiently small initialization scale $\epsilon = \tilde{\Theta}(\alpha_0^{8K})$, leads to a solution 360 $\boldsymbol{\theta}(t), t \geq 0$ such that: for some $t^* = \tilde{\mathcal{O}}\left(\log \frac{1}{\alpha_0}\right)$ and $T^* = \tilde{\Theta}\left(\log \frac{1}{\alpha_0}\right) + \tilde{\Omega}\left(\frac{1}{\alpha_0^{\min\{p-2,2\}}}\right)$ with 361 362 $[t^*, T^*] \neq \emptyset$, we have $\mathcal{L}(\boldsymbol{\theta}(t)) = \tilde{\mathcal{O}}(\alpha_0^4), \forall t \in [t^*, T^*]$ and 363 364

36

$$\sup_{\in [t^*,T^*]} \sup_{\boldsymbol{x} \in \mathbb{S}^{D-1}} \left| f^{(p)}(\boldsymbol{x};\boldsymbol{\theta}(t)) - F^{(p)}(\boldsymbol{x}) \right| \le \tilde{\mathcal{O}}\left(\alpha_0^2\right) \,. \tag{12}$$

368 $\hat{\Omega}, \tilde{o}, \tilde{O}$ hide logarithmic factor $\log \frac{K}{\delta}$ and constant factors that depend on p (in the worst case, 2^p). 369 We organize the subsequent discussions as follows: First, we state several remarks on understanding 370 our main result and comparing it with prior work; Then we move to a more technical discussion on 371 its proof sketch in Section 3.2.3; Lastly, we state in Section 3.2.4 several technical limitations of our 372 results and suggesting improvement in future research.

Nearly optimal robust classifier via GF The major implication of Theorem 3 is that one can find a nearly optimal ℓ_2 -robust classifier by GF without adversarial examples. When the intra-cluster variance α^2 is small, along the GF trajectory there exists a $f^{(p)}(\cdot; \boldsymbol{\theta}(t))$ that is $\tilde{\mathcal{O}}(\alpha_0^2)$ close to a nearly optimal ℓ_2 -robust classifier $F^{(p)}$, and such proximity to $F^{(p)}$ implies the same level of ℓ_2 -robustness, as shown in Proposition 1. We can immediately conclude that $f^{(p)}$ is also nearly optimal ℓ_2 -robust: **Corollary 1** (Nearly optimal ℓ_2 -robustness). Given any $f^{(p)}(\cdot; \boldsymbol{\theta}(t))$ obtained at $t \in [t^*, T^*]$ from Theorem 3, it can defend against adversarial attacks of radius $\frac{\sqrt{2}}{2} - \tilde{\mathcal{O}}(\alpha_0^{\frac{2}{p}})$ with probability $1 - \tilde{\mathcal{O}}(\alpha_0^{-2(1-\frac{2}{p})})$ over a new sample $(\boldsymbol{x}, \boldsymbol{y}) \sim \mathcal{D}_{X,Y}$, thus nearly optimal ℓ_2 -robust for $\mathcal{D}_{X,Y}$.

Comparison with prior work: robust classifier for orthonormal clusters The study of finding a 384 robust classifier for clusters with orthogonal cluster centers is initiated by Frei et al. (2023), where they theoretically show that any classifier obtained by gradient descent on a two-layer ReLU network is susceptible to an adversarial attack of ℓ_2 -radius $\mathcal{O}(\frac{1}{\sqrt{K}})$, despite that one can easily construct 386 387 a ReLU network that is robust to attacks of radius $\Theta(1)^2$. Then Min & Vidal (2024) explain this non-robustness issue of ReLU from a neural alignment perspective (Maennel et al., 2018; Boursier & 389 Flammarion, 2024), and propose pReLU to replace ReLU activation. They state as a conjecture that 390 training pReLU network $f^{(p)}(\cdot; \theta)$ under samples from $\mathcal{D}_{X,Y}$ leads to a classifier that is close to $F^{(p)}$ 391 when intra-cluster variance α^2 is small, and provide empirical validation to their conjecture. Our 392 work takes one step further to theoretically prove the convergence of pReLU towards $F^{(p)}$ under GF, 393 and also show that the achieved ℓ_2 -robustness is nearly optimal. Also, we believe a small initialization 394 is critical for finding a robust classifier as our data is approximately low-dimension thus adversarial examples exist if the initialization scale is large Melamed et al. (2024).

396

397 Comparison with prior work: GF on the two-layer network with small initialization Over 398 the past year, gradient descent/flow with small initialization has been studied for both linear net-399 works Gidel et al. (2019); Stöger & Soltanolkotabi (2021) and nonlinear networks Maennel et al. 400 (2018); Phuong & Lampert (2021); Boursier et al. (2022); Kumar & Haupt (2024); Chistikov et al. 401 (2023); Wang & Ma (2023); Min et al. (2024); Tsoy & Konstantinov (2024), to understand the implicit bias of gradient descent algorithms towards structurally simple networks. Our analysis follows 402 this line of work, as we will explain in Section 3.2.3 in detail, and also advances by considering a 403 more complicated dataset. Specifically, the GF on two-layer ReLU networks has been studied for 404 orthogonally separable data Phuong & Lampert (2021); Min et al. (2024); Chistikov et al. (2023), that 405 is, data with the same (different) label has positive (negative) correlation, for mutually orthogonal 406 data Boursier et al. (2022), and for positively correlated data (but only with two data points) Wang & 407 Ma (2023). Our data assumption is closest to mutually orthogonal data Boursier et al. (2022) (if we 408 set $\alpha = 0$), but considers a non-zero intra-cluster variance, which has not been studied in any of the 409 aforementioned work.

410 411

412

416

3.2.3 PROOF SKETCH

For simplicity, we consider the case $\alpha^2 = \alpha_0^2$ and use α^2 throughout this section. The discussion is conditioned on a good event (happens with probability at least $1 - \delta$) when samples are well concentrated around their respective cluster centers.

417 **Overall proof** Our proof in spirit is close to that of Boursier et al. (2022), with a two-phase 418 analysis of GF dynamics focusing on different quantities. Specifically, at the initial phase, called *alignment phase*, one studied the dynamics of the neuron direction $\frac{w_j}{\|w_j\|}$ through cosine angles 419 420 between w_i and cluster center μ_k , where one show, for all k and $j \in \mathcal{N}_k$, that $c_{kj} := \cos(\mu_k, w_j)$ 421 monotonically increases until it reaches $1 - \tilde{\mathcal{O}}(\alpha^2)$, that is, as we mentioned earlier, neurons 422 initialized within \mathcal{R}_k converge in angle to the corresponding μ_k . Then in the second *convergence* 423 phase, we show that all c_{kj} can probably stay above $1 - \tilde{\mathcal{O}}(\alpha^2)$ until T^* , and in the meantime, the norm of the neurons (measured by $\sum_{j \in \mathcal{N}_k} \| \boldsymbol{w}_j \|^2$ for each k) monotonically grow until reaches 424 425 $1 \pm \tilde{\mathcal{O}}(\alpha^2)$ before t^* . Moreover, the norm of the neurons initialized in the void region stays small: 426 $\sum_{j \in \mathcal{N}_c} \|\boldsymbol{w}_j(t)\| = \tilde{o}(\alpha^2). \text{ These three conditions } c_{kj} \ge 1 - \tilde{\mathcal{O}}(\alpha^2), \left|1 - \sum_{j \in \mathcal{N}_k} \|\boldsymbol{w}_j\|^2\right| \le \tilde{\mathcal{O}}(\alpha^2) \text{ and } \sum_{j \in \mathcal{N}_c} \|\boldsymbol{w}_j\| = \tilde{o}(\alpha^2) \text{ together imply the desired bound between } f^{(p)} \text{ and } F^{(p)}. \text{ We refer the } f^{(p)} \text{ and } F^{(p)} \text{ together imply the desired bound between } f^{(p)} \text{ and } F^{(p)}. \text{ We refer the } f^{(p)} \text{ together imply the desired bound between } f^{(p)} \text{ and } F^{(p)} \text{ together imply the desired bound between } f^{(p)} \text{ together imply } f^{(p)} \text{ to$ 427 428 readers to Figure 4 for an illustration of these phases. 429

² Frei et al. (2023) considers data sampled from $\mathcal{N}(\sqrt{D}\boldsymbol{\mu}_k, \alpha^2 \boldsymbol{I}), 1 \leq k \leq K$, thus their results should be rescaled by $\frac{1}{\sqrt{D}}$ when applied to $\mathcal{D}_{X,Y}$.

443 444

445

446

447

448

449 450 451

477

478



Figure 4: Important quantities (alignment and weight norms) and their dynamics long GF trajectory

Alignment phase (See Appendix E) During the alignment phase (the time interval between 0 and some $\tilde{O}(\log \frac{1}{\epsilon})$ time, where ϵ is the initialization scale). The norms of the weights stays $\tilde{O}(\epsilon)$ small, which follows a similar proof in Boursier et al. (2022); Min et al. (2024). The small norm bound on weights, together with the positive non-degeneracy gap assumption, allows the following characterization of the alignment for $1 \le k \le K, j \in \mathcal{N}_k$:

$$\frac{d}{dt}c_{kj} \ge Cpc_{kj}^{p-1}(1-c_{kj}^2) + \tilde{\mathcal{O}}\left(\frac{\alpha}{\sqrt{N}} + \frac{\alpha}{\sqrt{D}}\right) + \tilde{\mathcal{O}}\left(\alpha^2 + \frac{\alpha}{\sqrt{D}}\right) + \tilde{\mathcal{O}}\left(\epsilon\right), \quad (13)$$

for some constant C > 0. Consider the case when $\alpha = 0$, and $\epsilon \to 0$, for $j \in \mathcal{N}_k$, the dynamics $\frac{d}{dt}c_{kj} \ge Cpc_{kj}^{p-1}(1-c_{kj}^2)$ characterize the nominal effect of cluster centers $\mu_k, 1 \le k \le K$ on neuron direction $\frac{w_j}{\|w_j\|}$: each cluster centers is either attracting or repelling $\frac{w_j}{\|w_j\|}$, depending on whether their label matches the sign of v_j , and the aggregate effect is pushing $\frac{w_j}{\|w_j\|}$ towards μ_k , the closest cluster center to w_j in angle at initialization. We call k-th cluster the target cluster for w_j .

458 The rest of the terms are considered perturbations due to noisy samples around cluster centers and a 459 non-zero initialization scale: The first $\tilde{O}\left(\frac{\alpha}{\sqrt{N}} + \frac{\alpha}{\sqrt{D}}\right)$ term is due to the noisy samples from (the 460 target) k-th cluster. Since we have a $\Delta = \Theta(1)$ non-degeneracy gap, w_j has a positive inner product 461 with every sampled data within the k-th cluster, then one can utilize concentration results to bound 462 the effect of noise. The second $\tilde{\mathcal{O}}\left(\alpha^2 + \frac{\alpha}{\sqrt{D}}\right)$ term is due to the noisy samples from other non-target 463 464 clusters. Unfortunately, we have no control over how many of them have positive inner products with 465 w_i , thus a worse bound $O(\alpha^2)$ is derived. Lastly $O(\epsilon)$ is due to an ϵ -small weight norm because the 466 nominal effect is derived when weight norms are all zero. With $N = \tilde{\Omega}(\alpha^{-2})$ samples, $D = \tilde{\Omega}(\alpha^{-2})$ 467 dimension, and small ϵ , the dominant terms become $\mathcal{O}(\alpha^2)$, allowing us to prove the following: 468

Proposition 2 (Alignment in pReLU network). *Given the same assumptions as in Theorem 3 and consider the same GF solution* $\theta(t), t \ge 0$. *There exist some* $t_1 = \mathcal{O}\left(\log \frac{1}{\alpha}\right)$ *and* $t_2 = \mathcal{O}\left(\log \frac{1}{\epsilon}\right)$ *such that* $\forall k$ *and* $\forall j \in \mathcal{N}_k$, $\cos(\mu_k, w_j(t)) \ge 1 - \tilde{\mathcal{O}}(\alpha^2), \forall t \in [t_1, t_2]$.

We explicitly state the result during the alignment phase in Proposition 2 to highlight the difference between its described alignment for pReLU network (p > 2) to that of Boursier et al. (2022) for ReLU networks, where neurons are aligned with class average $\mu_{+} = \sum_{1 \le k \le K_1} \mu_k$ and $\mu_{-} = \sum_{K_1+1 \le k \le K} \mu_k$ instead of cluster centers.

Convergence phase (See Appendix F) During the convergence phase, the weight norm grows and exceeds ϵ -level, as suggested by the following dynamics:

$$\frac{d}{dt}\sum_{j\in\mathcal{N}_k}\|\boldsymbol{w}_j\|^2 = \left(1-\sum_{j\in\mathcal{N}_k}\|\boldsymbol{w}_j\|^2 + \tilde{\mathcal{O}}\left(\frac{\alpha}{\sqrt{N}}\right) + \tilde{\mathcal{O}}\left(\alpha^2\right) + \tilde{\mathcal{O}}\left(\alpha^p\right)\right)\sum_{j\in\mathcal{N}_k}\|\boldsymbol{w}_j\|^2, \quad (14)$$

which holds whenever $c_{kj} \ge 1 - \tilde{\mathcal{O}}(\alpha^2)$, $\forall k, j \in \mathcal{N}_k$. The nominal dynamics $\frac{d}{dt} \sum_{j \in \mathcal{N}_k} \|\boldsymbol{w}_j\|^2 = (1 - \sum_{j \in \mathcal{N}_k} \|\boldsymbol{w}_j\|^2) \sum_{j \in \mathcal{N}_k} \|\boldsymbol{w}_j\|^2$ describes the weight growth if $c_{kj} = 1, \forall k, j \in \mathcal{N}_k$ and $\alpha = 0$. Following nominal dynamics, $\sum_{j \in \mathcal{N}_k} \|\boldsymbol{w}_j\|^2$ converges to 1 for every k, minimizes the ℓ_2 -loss.

486 The rest of the terms are considered perturbations due to noisy samples around cluster centers and 487 the fact that alignment c_{kj} are only close to 1. The first $\tilde{\mathcal{O}}\left(\frac{\alpha}{\sqrt{N}}\right)$ is due to the noisy sample from 488 the target k-th cluster, the second $\tilde{\mathcal{O}}(\alpha^2)$ term is from imperfect alignment $c_{kj} \ge 1 - \tilde{\mathcal{O}}(\alpha^2)$, and 489 490 the last $\hat{\mathcal{O}}(\alpha^p)$ term is from the noisy sample from the non-target clusters (Notice that now w_i s are 491 almost orthogonal to non-target clusters, thus the effect of non-target clusters is smaller than during alignment phase). With $N = \tilde{\Omega}(\alpha^{-2})$ samples, the dominant terms become $\tilde{\mathcal{O}}(\alpha^{2})$, allowing us to 492 show that $\sum_{j \in \mathcal{N}_{k}} \|\boldsymbol{w}_{j}\|^{2}$ converges to $1 \pm \tilde{\mathcal{O}}(\alpha^{2})$ within t^{*} time. 493 494

The only missing piece is that this argument requires $c_{kj} \ge 1 - \tilde{\mathcal{O}}(\alpha^2)$, $\forall k, j \in \mathcal{N}_k$ but one no longer has (13) after $\tilde{\Theta}(\log \frac{1}{\epsilon})$ when weight norm starts to grow to $\tilde{\Theta}(1)$ -level. Nonetheless, once the alignment c_{kj} is $1 - \tilde{\mathcal{O}}(\alpha^2)$, it is hard to drop below this level as it relies on the attraction from non-target clusters but they are now near orthogonal to the neurons. Indeed, during the convergence phase, we can show that $\frac{d}{dt}c_{kj} \ge -\tilde{\mathcal{O}}(\alpha^{\min\{p,4\}})$, by which we show c_{kj} can stay at $1 - \tilde{\mathcal{O}}(\alpha^2)$ level until T^* time. Since $T^* \ge t^*$ for small α , our analysis of the weight norm growth is valid.

501 502

504

505

3.2.4 TECHNICAL LIMITATIONS OF CURRENT RESULTS

We conclude by discussing several technical limitations of our current results and potential avenues to address them. These limitations are listed in an order that the most challenging ones are stated first.

506 **Requirement on the initialization** The initialization requires a non-degeneracy gap $\Delta = \Theta(1)$, 507 which generally cannot be achieved by random initialization: the cosines between neurons and cluster 508 centers are $\mathcal{O}(\frac{1}{\sqrt{D}})$ with high probability. Given that $D = \tilde{\Omega}(\alpha^{-2})$, the actually non-degeneracy 509 gap of a random initialization is $\tilde{\mathcal{O}}(\alpha)$. We have discussed this issue when we define non-degenerate 510 initialization in Section 3.2.1: When neurons are initialized close to the boundary between a Voronoi 511 region \mathcal{R}_k and another region \mathcal{R}_l (or the void region \mathcal{R}°), whether they align with μ_k or with μ_l (or 512 get further into void region) depends on the actually sampled points in the dataset. In this regard, 513 when weights are randomly initialized, there is a "burn-in" phase during which neurons "choose" 514 their target clusters depending on the samples, then once they get away from the boundary of these 515 Voronoi regions with $\Delta = \Theta(1)$ gap, we can characterize the GF dynamics afterward by Theorem 3.

516 517

Upper bound on N Regarding our requirement $\tilde{\Omega}(\alpha_0^{-2}) \leq N \leq \tilde{o}(\exp(\alpha_0^{-2}))$, we have discussed 518 the lower bound $N \ge \tilde{\Omega}(\alpha_0^{-2})$ in Section 3.2.3. In fact, one can remove this lower bound and get a 519 final bound $\tilde{\mathcal{O}}(\frac{\alpha_0}{\sqrt{N}})$ in Theorem 3. The upper bound $N \leq \tilde{o}(\exp(\alpha_0^{-2}))$ may seem puzzling. This 520 issue originates from ReLU nonlinearity: a data point must activate a neuron by having a positive inner 521 product. Our analysis requires that a neuron w_i is activated by every data point from its target cluster, 522 which is translated into two conditions: 1) $\Theta(1)$ non-degeneracy gap; and 2) $\sqrt{\log N}\alpha_0 = \tilde{o}(1)$. 523 Here $\sqrt{\log N}\alpha_0$ is essentially the radius of a ℓ_2 -ball centered at a cluster center that can contain all 524 the sampled points from that cluster with high probability. Without these conditions, there will be outliers in sampled points, which must be handled with extra analysis. We believe this is possible because those outliers will be rare and thus may have a negligible effect on the dynamics.

527

528 **Extension to classification losses** Our results for the alignment phase directly apply to classification 529 losses: The choice of the loss $\ell(y, \hat{y})$ only affects the alignment dynamics through $\nabla_{\hat{u}} \ell(y, \hat{y})|_{\hat{u}=0}$, and 530 this quantity is same (may up to a constant scaling) regardless of whether ℓ is exponential, logistic, 531 or ℓ_2 . However, the analysis of convergence phase critically depends on ℓ : Recall that in Section 3.2.3 we show that the nominal weight norm dynamics are $\dot{z} = (1-z)z, z = \sum_{j \in \mathcal{N}_k} \|w_j\|^2$ for 532 ℓ_2 loss. For exponential loss, the nominal dynamics become $\dot{z} = \exp(-z)z$, whose closed-form 533 solution is not available. A better characterization of the solution to the nominal dynamics of the type 534 $\dot{z} = \exp(-z)z$ in future research naturally leads to an extension of Theorem 3 to classification losses. 535

536

Analysis until finite time T^* Our focus is on the distance between $f^{(p)}(\cdot; \theta(t))$ and $F^{(p)}$, thus we restrict to the time interval $[0, T^*]$ when we have explicit control of all relevant quantities (alignment, weight norms, etc.). To show convergence towards a minimizer of the loss after T^* , we believe applying the results in Chatterjee (2022) suffices, following the approach in Boursier et al. (2022).

540 REFERENCES

- Anish Athalye, Nicholas Carlini, and David Wagner. Obfuscated gradients give a false sense of
 security: Circumventing defenses to adversarial examples. In *International conference on machine learning*, pp. 274–283. PMLR, 2018.
- Etienne Boursier and Nicolas Flammarion. Early alignment in two-layer networks training is a two-edged sword. *arXiv preprint arXiv:2401.10791*, 2024.
- Etienne Boursier, Loucas Pullaud-Vivien, and Nicolas Flammarion. Gradient flow dynamics of
 shallow relu networks for square loss and orthogonal inputs. In *Advances in Neural Information Processing Systems*, volume 35, pp. 20105–20118, 2022.
- Nicholas Carlini, Anish Athalye, Nicolas Papernot, Wieland Brendel, Jonas Rauber, Dimitris Tsipras, Ian Goodfellow, Aleksander Madry, and Alexey Kurakin. On evaluating adversarial robustness. *arXiv preprint arXiv:1902.06705*, 2019.
- Patrick Chao, Alexander Robey, Edgar Dobriban, Hamed Hassani, George J Pappas, and Eric Wong.
 Jailbreaking black box large language models in twenty queries. *arXiv preprint arXiv:2310.08419*, 2023.
- Sourav Chatterjee. Convergence of gradient descent for deep neural networks. *arXiv preprint arXiv:2203.16462*, 2022.
- Dmitry Chistikov, Matthias Englert, and Ranko Lazic. Learning a neuron by a shallow reLU network: Dynamics and implicit bias for correlated inputs. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
- Jeremy Cohen, Elan Rosenfeld, and Zico Kolter. Certified adversarial robustness via randomized smoothing. In *international conference on machine learning*, pp. 1310–1320. PMLR, 2019.
- Francesco Croce and Matthias Hein. Reliable evaluation of adversarial robustness with an ensemble
 of diverse parameter-free attacks. In *ICML*, 2020.
- Will Cukierski. Dogs vs. cats. https://kaggle.com/competitions/dogs-vs-cats, 2013. Kaggle.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In 2009 IEEE Conference on Computer Vision and Pattern Recognition,
 pp. 248–255, 2009. doi: 10.1109/CVPR.2009.5206848.
- Mahyar Fazlyab, Manfred Morari, and George J Pappas. Safety verification and robustness analysis of neural networks via quadratic constraints and semidefinite programming. *IEEE Transactions on Automatic Control*, 67(1):1–15, 2020.
- Spencer Frei, Gal Vardi, Peter Bartlett, and Nathan Srebro. The double-edged sword of implicit
 bias: Generalization vs. robustness in reLU networks. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
- Tomer Galanti, András György, and Marcus Hutter. On the role of neural collapse in transfer learning. In *International Conference on Learning Representations*, 2021.
- Gauthier Gidel, Francis Bach, and Simon Lacoste-Julien. Implicit regularization of discrete gradient
 dynamics in linear neural networks. In *Advances in Neural Information Processing Systems*,
 volume 32, pp. 3202–3211. Curran Associates, Inc., 2019.
- Sven Gowal, Chongli Qin, Jonathan Uesato, Timothy Mann, and Pushmeet Kohli. Uncovering the limits of adversarial training against norm-bounded adversarial examples. *arXiv preprint arXiv:2010.03593*, 2020.
- 592 Chuan Guo, Mayank Rana, Moustapha Cisse, and Laurens van der Maaten. Countering adversarial
 593 images using input transformations. In *International Conference on Learning Representations*, 2018.

604

610

622

- 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*, pp. 770–778, 2016.
- Jiabao Ji, Bairu Hou, Alexander Robey, George J Pappas, Hamed Hassani, Yang Zhang, Eric
 Wong, and Shiyu Chang. Defending large language models against jailbreak attacks via semantic
 smoothing. *arXiv preprint arXiv:2402.16192*, 2024.
- Kaleab A Kinfu and René Vidal. Analysis and extensions of adversarial training for video classifica tion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*,
 pp. 3416–3425, 2022.
- Akshay Kumar and Jarvis Haupt. Directional convergence near small initializations and saddles in two-homogeneous neural networks. *Transactions on Machine Learning Research*, 2024. ISSN 2835-8856.
- Alexander Levine and Soheil Feizi. (de) randomized smoothing for certifiable defense against patch
 attacks. *Advances in Neural Information Processing Systems*, 33:6465–6475, 2020.
- Jiangyuan Li, Thanh V Nguyen, Chinmay Hegde, and Raymond K. W. Wong. Implicit sparse regularization: The impact of depth and early stopping. In A. Beygelzimer, Y. Dauphin, P. Liang, and J. Wortman Vaughan (eds.), Advances in Neural Information Processing Systems, 2021. URL https://openreview.net/forum?id=QM8oG0bz1o.
- Yuanzhi Li, Tengyu Ma, and Hongyang Zhang. Algorithmic regularization in over-parameterized matrix sensing and neural networks with quadratic activations. In Sébastien Bubeck, Vianney Perchet, and Philippe Rigollet (eds.), *Proceedings of the 31st Conference On Learning Theory*, volume 75 of *Proceedings of Machine Learning Research*, pp. 2–47. PMLR, 06–09 Jul 2018. URL https://proceedings.mlr.press/v75/li18a.html.
- Hartmut Maennel, Olivier Bousquet, and Sylvain Gelly. Gradient descent quantizes relu network
 features. *arXiv preprint arXiv:1803.08367*, 2018.
- Odelia Melamed, Gilad Yehudai, and Gal Vardi. Adversarial examples exist in two-layer relu
 networks for low dimensional linear subspaces. Advances in Neural Information Processing
 Systems, 36, 2024.
- Hancheng Min and René Vidal. Can implicit bias imply adversarial robustness? In *Proceedings of the 41st International Conference on Machine Learning*, volume 235 of *Proceedings of Machine Learning Research*, pp. 35687–35718. PMLR, 21–27 Jul 2024.
- Hancheng Min, Enrique Mallada, and René Vidal. Early neuron alignment in two-layer relu networks
 with small initialization. In *International Conference on Learning Representations*, pp. 1–8, 5 2024.
- Edward Moroshko, Blake E Woodworth, Suriya Gunasekar, Jason D Lee, Nati Srebro, and Daniel
 Soudry. Implicit bias in deep linear classification: Initialization scale vs training accuracy. Advances
 in Neural Information Processing Systems, 33, 2020.
- Andrew Y. Ng, Michael I. Jordan, and Yair Weiss. On spectral clustering: analysis and an algorithm. In *Proceedings of the 14th International Conference on Neural Information Processing Systems: Natural and Synthetic*, NIPS'01, pp. 849–856, Cambridge, MA, USA, 2001. MIT Press.
- Ambar Pal, Jeremias Sulam, and Rene Vidal. Adversarial examples might be avoidable: The role of
 data concentration in adversarial robustness. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
- Ambar Pal, René Vidal, and Jeremias Sulam. Certified robustness against sparse adversarial perturbations via data localization. *arXiv preprint arXiv:2405.14176*, 2024.
- Nicolas Papernot, Patrick McDaniel, Xi Wu, Somesh Jha, and Ananthram Swami. Distillation as a defense to adversarial perturbations against deep neural networks. In 2016 IEEE symposium on security and privacy (SP), pp. 582–597. IEEE, 2016.

648 649 650	Vardan Papyan, XY Han, and David L Donoho. Prevalence of neural collapse during the terminal phase of deep learning training. <i>Proceedings of the National Academy of Sciences</i> , 117(40): 24652–24663, 2020.
651 652 653	Mary Phuong and Christoph H Lampert. The inductive bias of relu networks on orthogonally separable data. In <i>International Conference on Learning Representations</i> , 2021.
654 655	Alexander Robey, Eric Wong, Hamed Hassani, and George J Pappas. Smoothllm: Defending large language models against jailbreaking attacks. <i>arXiv preprint arXiv:2310.03684</i> , 2023.
657 658 659	Ali Shafahi, Mahyar Najibi, Mohammad Amin Ghiasi, Zheng Xu, John Dickerson, Christoph Studer, Larry S Davis, Gavin Taylor, and Tom Goldstein. Adversarial training for free! <i>Advances in</i> <i>Neural Information Processing Systems</i> , 32, 2019.
660 661 662	Rusheb Shah, Soroush Pour, Arush Tagade, Stephen Casper, Javier Rando, et al. Scalable and transferable black-box jailbreaks for language models via persona modulation. <i>arXiv preprint arXiv:2311.03348</i> , 2023.
663 664 665 666	Dominik Stöger and Mahdi Soltanolkotabi. Small random initialization is akin to spectral learning: Optimization and generalization guarantees for overparameterized low-rank matrix reconstruction. <i>Advances in Neural Information Processing Systems</i> , 34, 2021.
667 668	Jeremias Sulam, Ramchandran Muthukumar, and Raman Arora. Adversarial robustness of supervised sparse coding. <i>Advances in neural information processing systems</i> , 33:2110–2121, 2020.
669 670 671 672	Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and Rob Fergus. Intriguing properties of neural networks. In <i>2nd International Conference on Learning Representations</i> , 2014.
673 674 675	Florian Tramèr, Alexey Kurakin, Nicolas Papernot, Ian Goodfellow, Dan Boneh, and Patrick Mc- Daniel. Ensemble adversarial training: Attacks and defenses. In <i>International Conference on Learning Representations</i> , 2018.
676 677 678	Nikita Tsoy and Nikola Konstantinov. Simplicity bias of two-layer networks beyond linearly separable data. <i>arXiv preprint arXiv:2405.17299</i> , 2024.
679 680 681	Mingze Wang and Chao Ma. Understanding multi-phase optimization dynamics and rich nonlinear behaviors of reLU networks. In <i>Thirty-seventh Conference on Neural Information Processing Systems</i> , 2023.
682 683	Eric Wong, Leslie Rice, and J Zico Kolter. Fast is better than free: Revisiting adversarial training. In <i>International Conference on Learning Representations</i> , 2019.
685 686	Dongxian Wu, Shu-Tao Xia, and Yisen Wang. Adversarial weight perturbation helps robust general- ization. <i>Advances in neural information processing systems</i> , 33:2958–2969, 2020.
687 688 689	Greg Yang, Tony Duan, J Edward Hu, Hadi Salman, Ilya Razenshteyn, and Jerry Li. Randomized smoothing of all shapes and sizes. In <i>International Conference on Machine Learning</i> , pp. 10693–10705. PMLR, 2020.
690 691 692 693	Huan Zhang, Tsui-Wei Weng, Pin-Yu Chen, Cho-Jui Hsieh, and Luca Daniel. Efficient neural network robustness certification with general activation functions. <i>Advances in neural information processing systems</i> , 31, 2018.
694 695 696 697 698	Difan Zou, Spencer Frei, and Quanquan Gu. Provable robustness of adversarial training for learning halfspaces with noise. In <i>International Conference on Machine Learning</i> , pp. 13002–13011. PMLR, 2021.
699 700 701	

NUMERICAL EXPERIMENTS A

A.1 ADDITIONAL EXPERIMENTS ON LEARNING ROBUST CLASSIFIER FOR DATA FROM **ORTHONORMAL CLUSTERS**

In this section, we provide additional experiments to that in Figure 2, highlighting the importance of parametrization of function space, and the hyperparameters of training algorithm in determining whether one can succeed in obtaining robust classifier for data from orthonormal clusters.



Figure 5: Given sampled data from (1) with 12 positive clusters and 8 negative clusters (D = 2000), 723 gradient descent (SGD, small initialization) on (bias-free, width-200) two-layer network with regular 724 polynomial ReLU activation of degree 3 fails to find a robust classifier. Moreover, if one increases the 725 variance of the random initialization, both regular polynomial ReLU network and pReLU network 726 can not find a robust classifier. All networks here are trained for a sufficient amount of epochs until 727 they achieve perfect training accuracy on a synthesis dataset of our orthonormal cluster model of size 728 20000.729

Regular polynomial ReLU networks In this experiment, we consider both the regular polynomial 730 ReLU networks to pReLU networks. In particular, recall that the regular polynomial ReLU networks are defined as: 732

733 734

755

731

702

703 704

705

706

708

709

$$g(oldsymbol{x}; ilde{oldsymbol{ heta}}) = \sum_{j=1}^h v_j \sigma^p(\langle oldsymbol{x}, oldsymbol{w}_j
angle), \qquad (ilde{oldsymbol{ heta}} := \{oldsymbol{w}_j, v_j\}_{j=1}^h)\,.$$

735 (Two-layer Networks with Polynomial ReLU activation with degree *p*) 736 We note its difference with pReLU networks: regular polynomial ReLU networks do not have a weight normalization at the first layer. Nonetheless, when p is fixed, it is easy to verify that the 737 function/hypothesis spaces induced by pReLU networks and regular polynomial ReLU networks 738 are the same: any function $f^{(p)}(\boldsymbol{x};\boldsymbol{\theta})$ for some $\boldsymbol{\theta} = \{\boldsymbol{w}_j, v_j\}_{j=1}^h$ is equivalent to $g(\boldsymbol{x}; \tilde{\boldsymbol{\theta}})$ with 739 $\tilde{\boldsymbol{\theta}} = \{\boldsymbol{w}_j, \frac{v_j}{\|\boldsymbol{w}_i\|^{p-1}}\}_{j=1}^{h-3}.$ 740

741 Regular polynomial ReLU networks v.s. pReLU Although the induced function/hypothesis spaces 742 are the same, GD on regular polynomial ReLU networks and pReLU finds classifiers with different 743 levels of robustness. As one can see in Figure 5, with a small initialization (all weight entries are 744 randomly initialized as $\mathcal{N}(0, 1 \times 10^{-4})$), SGD on a pReLU network successfully finds a classifier that 745 is as robust as the Bayes classifier. However, SGD on a regular polynomial ReLU network fails to find 746 a robust classifier. This suggests that the way the function/hypothesis spaces are parametrized is also 747 important in determining the robustness of the networks trained by GD, as different parametrization 748 induces different implicit biases of GD in selecting the loss minimizer in the function space.

749 **Effect of initialization scale** Finally, when one uses a large initialization scale, where all weight 750 entries are randomly initialized as $\mathcal{N}(0, 0.25)$, even the GD on a pReLU network fails to find a 751 robust classifier. This is not surprising as the initialization scale also controls the implicit bias of 752 GD Moroshko et al. (2020), and many works Maennel et al. (2018); Stöger & Soltanolkotabi (2021); 753 Li et al. (2018; 2021) have theoretically shown the advantage of using a small initialization scale in 754 GD.

³The neurons with $\|\boldsymbol{w}_i\| = 0$ should be eliminated from the parameters for this argument to hold.

756 A.2 CAT V.S. DOG CLASSIFICATION VIA TRANSFER LEARNING

758 In this section, we solve the tasks of classifying cats and dogs (Cukierski, 2013) via transfer learning 759 using extracted features from a ResNet152 (He et al., 2016) trained on ImageNet (Deng et al., 2009). We conjecture that the extracted features of the dog (or cat) class may naturally have many 760 clusters: when the feature extractor is trained on ImageNet, dogs are further labeled by their breeds. 761 Thus the extracted features of dogs of the same breed should be sufficiently close, and features of 762 dogs of different breeds should be sufficiently far apart, based on the well-known neural collapse 763 phenomenon (Papyan et al., 2020; Galanti et al., 2021). If such a multi-cluster structure exists in the 764 extracted feature, then we expect training pReLU as a classification head can achieve better robust 765 accuracy compared to its ReLU counterpart. 766

The rest of the section is organized as follows: First, we show that the extracted features of cats v.s. dogs dataset exhibits a multi-cluster structure; Then we train pReLU networks with different choices of p as a classification head and compute the robust accuracy of these train networks with AutoAttack (Croce & Hein, 2020) on the extracted feature space.



771

772 773

774

775

776

777

778

779

781

782

783

784

785

786

787

788 789

790

791

792 793

794

796

797

798 799

800

801

802

803

804

805

806 807

Figure 6: Pairwise inner product of the features of 3000 images of **dogs** from cat v.s. dog dataset Cukierski (2013). The features are clustered into 9 clusters via spectral clustering.



Figure 8: Pairwise inner product of the features of 3000 images of **cats** from cat v.s. dog dataset Cukierski (2013). The features are clustered into 10 clusters via spectral clustering.



Figure 7: Average pairwise inner product between features from two clusters of **dogs** features. The features are clustered into 9 clusters; Each pixel $(i, j), 1 \le i, j \le 9$ represents the average inner product between features from cluster *i* and cluster *j*.



Figure 9: Average pairwise inner product between features from two clusters of **cats** features. The features are clustered into 10 clusters; Each pixel $(i, j), 1 \le i, j \le 10$ represents the average inner product between features from cluster *i* and cluster *j*.

Multi-cluster structure of extracted feature We first collect extracted features of the entire cat
 v.s. dog dataset (Cukierski, 2013), center these features by the global mean feature vector and then normalized all the features. Then we take a subset of the centered, normalized features (for the sake

of simplicity, we call the centered, normalized features as features) from the same class (cat or dog), do spectral clustering (Ng et al., 2001) on the features, then compute the inner product between the features. From Figure 6 and 7, we see even within the same (dog) class, the extracted features have a multi-cluster structure, and we conjecture that this is because when the feature extractor is trained on ImageNet, dogs are further labeled by their breeds. Interestingly, if we perform the same visualization for cat images, as in Figure 8 and 9, the multi-cluster structure still exists but with less prominent clusters; We conjecture that this is because ImageNet has much less cat classes than dog classes.

Training pReLU as classification head Now, with the extracted features of cat v.s. dog dataset, we train two-layer pReLU networks with different choices of p using Adam, following the same experiment settings in (Min & Vidal, 2024). After training, we compute the robust accuracy of the trained networks under adaptive adversarial ℓ_2 and ℓ_{∞} attacks (Croce & Hein, 2020). We observe that pReLU networks with larger p achieve better robust accuracy than ReLU networks (p = 1).



Figure 10: Cat and dog classification: Training and test accuracy v.s. training epochs

Figure 11: Robust accuracy of trained networks under ℓ_2 PGD attacks.

Figure 12: Robust accuracy of trained networks under ℓ_{∞} PGD attacks.

In summary, we show that in a transfer learning scenario, the multi-cluster structure arises due to the distinguishing power of the feature extractor trained on large datasets with finer labels, and we show that in this case, pReLU networks with larger p achieve better robustness compared to its ReLU counterpart. Admittedly, our current Theorems cannot fully explain the observed experimental results since the extracted features form clusters with large variances, and there are some correlations among these clusters, which does not follow our data assumption. Relaxing our data assumption to large variance, and allowing inter-cluster correlation is an import future research direction.



OPTIMAL ROBUST CLASSIFIER FOR ORTHONORMAL CLUSTERS В

In this section, we discuss the optimal robust classifier for orthonormal clusters. We first show that any measurable classifier can not defend against an adversarial attack of ℓ_2 radius $\frac{\sqrt{2}}{2}$, leading to a robust error of at least $\frac{\min\{K_1, K_2\}}{K}$. Then we consider the Bayes optimal classifier $f^*(x) = \arg \max_y \mathbb{P}(Y = y | x)$ and show that it is also optimally robust: it can defend against any adversarial attack of ℓ_2 radius $\frac{\sqrt{2}}{2} - o(1)$, as the dimension of the data D increases.

B.1 Maximum robustness against ℓ_2 adversarial attacks

We need the following lemma (we provide proof after proving Theorem 1)

Lemma 1. For any $n \times m$ matrix, let a be the number of rows that contain at least one nonpositive entry and b be the number of columns that contain at least one non-negative entry. Then $a+b \ge \min\{n,m\}.$

With Lemma 1, we are ready to prove Theorem 1.

Theorem 1 (Restated). Let $f : \mathbb{R}^D \to \mathbb{R}$ be any Lebesgue measurable function such that the random variable $\min_{\|\boldsymbol{d}\|\leq 1} \left[f\left(\boldsymbol{x} + \frac{\sqrt{2}}{2}\boldsymbol{d}\right) y \right]$ is also measurable. Given a sample $(\boldsymbol{x}, y) \sim \mathcal{D}_{X,Y}$, we have Г

$$\mathbb{P}_{(\boldsymbol{x},\boldsymbol{y})\sim\mathcal{D}_{\boldsymbol{X},\boldsymbol{Y}}}\left(\min_{\|\boldsymbol{d}\|\leq 1}\left[f\left(\boldsymbol{x}+\frac{\sqrt{2}}{2}\boldsymbol{d}\right)\boldsymbol{y}\right]\leq 0\right)\geq\frac{\min\{K_1,K_2\}}{K}.$$
(B.1)

Proof. We start with the following:

$$\mathbb{P}\left(\min_{\|\boldsymbol{d}\|\leq 1}\left[f\left(\boldsymbol{x}+\frac{\sqrt{2}}{2}\boldsymbol{d}\right)\boldsymbol{y}\right]\leq 0\right)=\sum_{k=1}^{K}\mathbb{P}\left(\min_{\|\boldsymbol{d}\|\leq 1}\left[f\left(\boldsymbol{x}+\frac{\sqrt{2}}{2}\boldsymbol{d}\right)\boldsymbol{y}\right]\leq 0\mid \boldsymbol{z}=\boldsymbol{k}\right)\mathbb{P}\left(\boldsymbol{z}=\boldsymbol{k}\right)$$
(B.2)

For $k \leq K_1$,

$$\mathbb{P}\left(\min_{\|\boldsymbol{d}\|\leq 1}\left[f\left(\boldsymbol{x}+\frac{\sqrt{2}}{2}\boldsymbol{d}\right)\boldsymbol{y}\right]\leq 0\mid \boldsymbol{z}=\boldsymbol{k}\right) = \mathbb{P}_{\boldsymbol{\varepsilon}}\left(\min_{\|\boldsymbol{d}\|\leq 1}\left[f\left(\boldsymbol{\mu}_{k}+\boldsymbol{\varepsilon}+\frac{\sqrt{2}}{2}\boldsymbol{d}\right)\right]\leq 0\right) \\ \geq \mathbb{P}_{\boldsymbol{\varepsilon}}\left(\min_{K_{1}+1\leq l\leq K}\left[f\left(\frac{\boldsymbol{\mu}_{k}+\boldsymbol{\mu}_{l}}{2}+\boldsymbol{\varepsilon}\right)\right]\leq 0\right).$$

The measurability of f ensures this lower bound exists. Similarly, we have for $K_1 + 1 \le k \le K$

$$\mathbb{P}\left(\min_{\|\boldsymbol{d}\|\leq 1}\left[f\left(\boldsymbol{x}+\frac{\sqrt{2}}{2}\boldsymbol{d}\right)\boldsymbol{y}\right]\leq 0\mid \boldsymbol{z}=\boldsymbol{k}\right)=\mathbb{P}_{\boldsymbol{\varepsilon}}\left(\min_{\|\boldsymbol{d}\|\leq 1}\left[-f\left(\boldsymbol{\mu}_{k}+\boldsymbol{\varepsilon}+\frac{\sqrt{2}}{2}\boldsymbol{d}\right)\right]\leq 0\right)$$
$$\geq \mathbb{P}_{\boldsymbol{\varepsilon}}\left(\min_{1\leq l\leq K_{1}}\left[-f\left(\frac{\boldsymbol{\mu}_{k}+\boldsymbol{\mu}_{l}}{2}+\boldsymbol{\varepsilon}\right)\right]\leq 0\right)$$
$$=\mathbb{P}_{\boldsymbol{\varepsilon}}\left(\max_{1\leq l\leq K_{1}}\left[f\left(\frac{\boldsymbol{\mu}_{k}+\boldsymbol{\mu}_{l}}{2}+\boldsymbol{\varepsilon}\right)\right]\geq 0\right).$$

Therefore,

$$\mathbb{P}\left(\min_{\|\boldsymbol{d}\|\leq 1} \left[f\left(\boldsymbol{x} + \frac{\sqrt{2}}{2}\boldsymbol{d}\right) \boldsymbol{y} \right] \leq 0 \right)$$
$$= \sum_{k=1}^{K} \mathbb{P}\left(\min_{\|\boldsymbol{d}\|\leq 1} \left[f\left(\boldsymbol{x} + \frac{\sqrt{2}}{2}\boldsymbol{d}\right) \boldsymbol{y} \right] \leq 0 \mid \boldsymbol{z} = \boldsymbol{k} \right) \mathbb{P}\left(\boldsymbol{z} = \boldsymbol{k}\right)$$
$$= \frac{1}{K} \left[\sum_{1\leq k\leq K_{1}} \mathbb{P}\left(\min_{\|\boldsymbol{d}\|\leq 1} \left[f\left(\boldsymbol{x} + \frac{\sqrt{2}}{2}\boldsymbol{d}\right) \boldsymbol{y} \right] \leq 0 \mid \boldsymbol{z} = \boldsymbol{k} \right)$$

$$\begin{aligned} +\sum_{K_{1}+1\leq k\leq K} \mathbb{P}\left(\min_{\|\boldsymbol{d}\|\leq 1}\left[f\left(\boldsymbol{x}+\frac{\sqrt{2}}{2}\boldsymbol{d}\right)\boldsymbol{y}\right]\leq 0\mid\boldsymbol{z}=\boldsymbol{k}\right)\right]\\ \geq \frac{1}{K}\left[\sum_{1\leq k\leq K_{1}} \mathbb{P}_{\boldsymbol{\varepsilon}}\left(\min_{K_{1}+1\leq l\leq K}\left[f\left(\frac{\boldsymbol{\mu}_{k}+\boldsymbol{\mu}_{l}}{2}+\boldsymbol{\varepsilon}\right)\right]\leq 0\right)\right.\\ &+\sum_{K_{1}+1\leq k\leq K} \mathbb{P}_{\boldsymbol{\varepsilon}}\left(\max_{1\leq l\leq K_{1}}\left[f\left(\frac{\boldsymbol{\mu}_{k}+\boldsymbol{\mu}_{l}}{2}+\boldsymbol{\varepsilon}\right)\right]\geq 0\right)\right]\\ &=\frac{1}{K}\left[\sum_{1\leq k\leq K_{1}}\int \mathbb{1}\left(\min_{K_{1}+1\leq l\leq K}\left[f\left(\frac{\boldsymbol{\mu}_{k}+\boldsymbol{\mu}_{l}}{2}+\boldsymbol{\varepsilon}\right)\right]\leq 0\right)\boldsymbol{p}(\boldsymbol{\varepsilon})\right.\\ &+\sum_{K_{1}+1\leq k\leq K}\int \mathbb{1}\left(\max_{1\leq l\leq K_{1}}\left[f\left(\frac{\boldsymbol{\mu}_{k}+\boldsymbol{\mu}_{l}}{2}+\boldsymbol{\varepsilon}\right)\right]\geq 0\right)\boldsymbol{p}(\boldsymbol{\varepsilon})\right]\\ &=\frac{1}{K}\int\left[\sum_{1\leq k\leq K_{1}}\mathbb{1}\left(\min_{K_{1}+1\leq l\leq K}\left[f\left(\frac{\boldsymbol{\mu}_{k}+\boldsymbol{\mu}_{l}}{2}+\boldsymbol{\varepsilon}\right)\right]\leq 0\right)\right) \tag{B.3}$$

$$\frac{1}{k \leq K_{1}} \left(\begin{array}{c} K_{1}+1 \leq l \leq K \end{array} \right) \left(\begin{array}{c} 2 \\ + \sum_{K_{1}+1 \leq k \leq K} \mathbb{1} \left(\max_{1 \leq l \leq K_{1}} \left[f\left(\frac{\mu_{k}+\mu_{l}}{2}+\varepsilon\right) \right] \geq 0 \right) \right] p(\varepsilon), \quad (B.4)$$

and if we define the $K_1 \times K_2$ matrix

$$M_f(\boldsymbol{\varepsilon}) := \left[f\left(\frac{\boldsymbol{\mu}_k + \boldsymbol{\mu}_l}{2} + \boldsymbol{\varepsilon}\right) \right]_{1 \le k \le K_1, \ K_1 + 1 \le l \le K}$$
(B.5)

and examine carefully enough, we notice that $\sum_{1 \le k \le K} \mathbb{1} \left(\min_{K_1+1 \le l \le K} \left[f \left(\frac{\mu_k + \mu_l}{2} + \varepsilon \right) \right] \le 0 \right)$ is the number of rows of $M_f(\varepsilon)$ that contains at least one non-positive entry and $\sum_{K_1+1 \le k \le K} \mathbb{1} \left(\max_{1 \le l \le K_1} \left[f \left(\frac{\mu_k + \mu_l}{2} + \varepsilon \right) \right] \ge 0 \right)$ is the number of columns of $M_f(\varepsilon)$ that contains at least one non-negative entry. By Lemma 1, we have

$$\mathbb{P}\left(\min_{\|\boldsymbol{d}\|\leq 1}\left[f\left(\boldsymbol{x}+\frac{\sqrt{2}}{2}\boldsymbol{d}\right)\boldsymbol{y}\right]\leq 0\right)\geq (\mathbf{B}.4)\geq \frac{1}{K}\int\min\{K_1,K_2\}p(\boldsymbol{\varepsilon})\,.$$

Therefore

$$\mathbb{P}\left(\min_{\|\boldsymbol{d}\|\leq 1}\left[f\left(\boldsymbol{x}+\frac{\sqrt{2}}{2}\boldsymbol{d}\right)\boldsymbol{y}\right]\leq 0\right)\geq \frac{\min\{K_1,K_2\}}{K}.$$
(B.6)

961 Proof of Lemma 1. We denote $C^*(n, m)$ the minimum value of a+b over all possible choice of $n \times m$ 962 matrices. It suffices to show $C^*(n, m) \ge \min\{n, m\}$ (The equality is obtained by an all-positive 963 matrix when $n \le m$ and an all-negative matrix otherwise), and we prove it by induction. 964 Denote the set of the se

For $n = 1, m = 1, C^*(n, m) = 1$. This is trivial. We need to show that if $C^*(n, m) = \min\{n, m\}$ holds for some n and m, then

• $C^*(n, m+1) = \min\{n, m+1\};$

• and
$$C^*(n+1,m) = \min\{n+1,m\}.$$

We shall prove these two cases:

972 **Case 1** $C^*(n,m) \ge \min\{n,m\} \Rightarrow C^*(n,m+1) \ge \min\{n,m+1\}$

Given an $n \times m$ matrix M and an agumented matrix $M' = [M \ v]$, we let a, b and a', b' be the row/column counts of our interest for M and M' respectively. Without loss of generality, we suppose the first a rows of M all contain at least one non-positive entry (and the rest do not, by definition of a). We know that $a + b \ge \min\{n, m\}$, and

$$a' = a + \sum_{i=a+1}^{n} \mathbb{1}(v_i \le 0), \qquad b' = b + \mathbb{1}(\max_i v_i \ge 0),$$
 (B.7)

which is

982 983 984

985 986

987 988

989 990

991

992

993

994 995

1002 1003

1004 1005

1007 1008

1014

1025

$$a' + b' = a + b + \sum_{i=a+1}^{n} \mathbb{1}(v_i \le 0) + \mathbb{1}(\max_i v_i \ge 0).$$
 (B.8)

There are two scenarios:

- 1. When a = n, we have $\sum_{i=a+1}^{n} \mathbb{1}(\boldsymbol{v}_i \leq 0) + \mathbb{1}(\max_i \boldsymbol{v}_i \geq 0) \geq 0$ 2. When a < n, we have $\sum_{i=a+1}^{n} \mathbb{1}(\boldsymbol{v}_i \leq 0) + \mathbb{1}(\max_i \boldsymbol{v}_i \geq 0) \geq 1$. Therefore, we find that
- $a'+b' \ge \min\{n+b, a+b+1\} \ge \min\{n, \min\{n, m\}+1\} = \min\{n, n+1, m+1\} = \min\{n, m+1\}.$ (B.9) This shows $C^*(n, m+1) \ge \min\{n, m+1\}.$

Case 2
$$C^*(n+1,m) \ge \min\{n+1,m\} \Rightarrow C^*(n+1,m) \ge \min\{n+1,m\}$$

Given an $n \times m$ matrix M and an agumented matrix $M' = \begin{bmatrix} M \\ v \end{bmatrix}$, we let a, b and a', b' be the row/column counts of our interest for M and M' respectively. Without loss of generality, we suppose

the first b columns of M all contain at least one non-negative entry (and the rest do not, by definition of b). We know that $a + b \ge \min\{n, m\}$, and

$$a' = a + \mathbb{1}(\min_{i} v_{i} \le 0), \qquad b' = b + \sum_{i=b+1}^{m} \mathbb{1}(v_{i} \ge 0),$$
 (B.10)

1006 which is

$$a' + b' = a + b + \sum_{i=b+1}^{m} \mathbb{1}(v_i \ge 0) + \mathbb{1}(\min_i v_i \le 0).$$
 (B.11)

There are two scenarios:

10111. When
$$b = m$$
, we have $\sum_{i=b+1}^{m} \mathbb{1}(v_i \ge 0) + \mathbb{1}(\min_i v_i \le 0) \ge 0$ 10122. When $b < m$, we have $\sum_{i=b+1}^{m} \mathbb{1}(v_i \ge 0) + \mathbb{1}(\min_i v_i \le 0) \ge 1$

1015 Therefore, we find that

$$\begin{array}{ll} \mbox{1016} & a'+b' \geq \min\{a+m,a+b+1\} \geq \min\{m,\min\{n,m\}+1\} = \min\{m,n+1,m+1\} = \min\{n+1,m\} & (B.12) \\ \mbox{1018} & \mbox{This shows } C^*(n+1,m) \geq \min\{n+1,m\}. & \mbox{\square} \end{array}$$

1020 B.2 BAYES OPTIMAL CLASSIFIER W.R.T. 0-1 LOSS 1021

1022 Our proof will use Hoeffding's inequality for high-dimensional Gaussian vectors 1023 Lemma 2 (Hoeffding inequality). For any unit vector $\mu \in \mathbb{S}^{D-1}$, we have

$$\mathbb{P}_{\boldsymbol{\varepsilon} \sim \mathcal{N}\left(0, \frac{\alpha^{2}}{D}\boldsymbol{I}\right)}\left(\left|\left\langle \boldsymbol{\mu}, \boldsymbol{\varepsilon} \right\rangle\right| > t\right) \leq 2\exp\left(-\frac{Dt^{2}}{2\alpha^{2}}\right).$$
(B.13)

And the concentration result of the norm of high-dimensional Gaussian vectors

Lemma 3. We have

$$\mathbb{P}_{\boldsymbol{\varepsilon} \sim \mathcal{N}\left(0,\frac{\alpha^2}{D}\boldsymbol{I}\right)}\left(\|\boldsymbol{\varepsilon}\| > t\right) \le 4\exp\left(-\frac{t^2}{8\alpha^2}\right),\tag{B.14}$$

Theorem 2 (Restated). The Bayes optimal classifier for label Y given observation x w.r.t. 0-1 loss is $\operatorname{sign}(f^*(\boldsymbol{x}))$, where $f^*(\boldsymbol{x}) = \sum_{k=1}^{K_1} \exp\left(\frac{D\langle \boldsymbol{x}, \boldsymbol{\mu}_k \rangle}{\alpha^2}\right) - \sum_{k=K_1+1}^{K} \exp\left(\frac{D\langle \boldsymbol{x}, \boldsymbol{\mu}_k \rangle}{\alpha^2}\right)$. Moreover, given a sample $(\boldsymbol{x}, y) \sim \mathcal{D}_{X,Y}$, we have, for any $\frac{2\sqrt{2\alpha^2 \log K}}{D} \leq \nu \leq \sqrt{2}$,

$$\mathbb{P}_{(\boldsymbol{x},\boldsymbol{y})\sim\mathcal{D}_{\boldsymbol{X},\boldsymbol{Y}}}\left(\min_{\|\boldsymbol{d}\|\leq 1}\left[f^*\left(\boldsymbol{x}+\frac{\sqrt{2}-\nu}{2}\boldsymbol{d}\right)\boldsymbol{y}\right]>0\right)\geq 1-2K\exp\left(-\frac{D\nu^2}{64\alpha^2}\right).$$
(B.15)

Proof. Bayes optimal classifier for $\mathcal{D}_{X,Y}$ The Bayes optimal classifier w.r.t. 0-1 loss is given by

$$f^{*}(\boldsymbol{x}) = \arg \max_{y} \mathbb{P} \left(Y = y \mid X = \boldsymbol{x} \right)$$

$$= \arg \max_{y} \sum_{k=1}^{K} \mathbb{P} \left(Y = y \mid Z = k, X = \boldsymbol{x} \right) \mathbb{P} \left(Z = k \mid X = \boldsymbol{x} \right)$$

$$= \begin{cases} 1, & \text{if } \sum_{k=1}^{K_{1}} \mathbb{P} \left(Z = k \mid X = \boldsymbol{x} \right) > \sum_{k=K_{1}+1}^{K} \mathbb{P} \left(Z = k \mid X = \boldsymbol{x} \right) \\ -1, & \text{o.w.} \end{cases}$$

$$= \operatorname{sign} \left(\sum_{k=1}^{K_{1}} \mathbb{P} \left(Z = k \mid X = \boldsymbol{x} \right) - \sum_{k=K_{1}+1}^{K} \mathbb{P} \left(Z = k \mid X = \boldsymbol{x} \right) \right). \quad (B.16)$$
we srule and a few derivations give:

1050 Bayes rule and a few derivations give:
1051
$$\mathbb{P}(X - m \mid Z)$$

$$\mathbb{P}(Z = k \mid X = x) = \frac{\mathbb{P}(X = x \mid Z = k) \mathbb{P}(Z = k)}{\sum_{l=1}^{K} \mathbb{P}(X = x \mid Z = l) \mathbb{P}(Z = l)} \\
= \frac{\exp\left(-\frac{D||x - \mu_k||^2}{2\alpha^2}\right)}{\sum_{l=1}^{K} \exp\left(-\frac{D||x - \mu_k||^2}{2\alpha^2}\right)} \\
= \frac{\exp\left(-\frac{D(||x||^2 - 2\langle x, \mu_k \rangle + ||\mu_k||^2)}{2\alpha^2}\right)}{\sum_{l=1}^{K} \exp\left(-\frac{D(||x||^2 - 2\langle x, \mu_l \rangle + ||\mu_k||^2)}{2\alpha^2}\right)} = \frac{\exp\left(\frac{D\langle x, \mu_k \rangle}{\alpha^2}\right)}{\sum_{l=1}^{K} \exp\left(\frac{D\langle x, \mu_l \rangle}{\alpha^2}\right)}. \quad (B.17)$$

Combining (B.16) and (B.17), we have

$$f^{*}(\boldsymbol{x}) = \operatorname{sign}\left(\sum_{k=1}^{K_{1}} \exp\left(\frac{D\langle \boldsymbol{x}, \boldsymbol{\mu}_{k} \rangle}{\alpha^{2}}\right) - \sum_{k=K_{1}+1}^{K} \exp\left(\frac{D\langle \boldsymbol{x}, \boldsymbol{\mu}_{k} \rangle}{\alpha^{2}}\right)\right).$$
(B.18)

k),

Robustness of f^* . We now proceed to show that f^* is robust near-optimally. Since

$$\begin{split} & \mathbb{P}\left(\min_{\|\boldsymbol{d}\|\leq 1}\left[f^*\left(\boldsymbol{x}+\frac{\sqrt{2}-\nu}{2}\boldsymbol{d}\right)y\right]\leq 0\right) \\ & = \sum_{k=1}^{K}\mathbb{P}\left(\min_{\|\boldsymbol{d}\|\leq 1}\left[f^*\left(\boldsymbol{x}+\frac{\sqrt{2}-\nu}{2}\boldsymbol{d}\right)y\right]\leq 0\bigg|\,z=k\right)\mathbb{P}\left(z=\frac{1}{2}\right) \end{split}$$

It suffices to show that $\forall 1 \leq k \leq K$

$$\mathbb{P}\left(\min_{\|\boldsymbol{d}\|\leq 1}\left[f^*\left(\boldsymbol{x}+\frac{\sqrt{2}-\nu}{2}\boldsymbol{d}\right)\boldsymbol{y}\right]\leq 0\left|\boldsymbol{z}=\boldsymbol{k}\right)\leq K\exp\left(-\frac{CD\nu^2}{16\alpha^2}\right).$$
(B.19)

When $k \leq K_1$, we have

$$\mathbb{P}\left(\min_{\|\boldsymbol{d}\|\leq 1}\left[f^*\left(\boldsymbol{x}+\frac{\sqrt{2}-\nu}{2}\boldsymbol{d}\right)y\right]\leq 0 \left|z=k\right)$$

$$\begin{split} &= \mathbb{P}_{\varepsilon} \left(\min_{\|d\| \leq 1} \left[f^{\varepsilon} \left(x + \frac{\sqrt{2} - \nu}{2} d \right) \right] \leq 0 \right) \\ &= \mathbb{P}_{\varepsilon} \left(\min_{\|d\| \leq 1} \left[\exp \left(\frac{D}{\alpha^{2}} \left(1 + \langle \mu_{k}, \varepsilon \rangle + \frac{\sqrt{2} - \nu}{2} \langle d, \mu_{k} \rangle \right) \right) \right) \\ &\quad + \sum_{i \neq k, 1 \leq l \leq K} \exp \left(\frac{D}{\alpha^{2}} \left(\langle \mu_{i}, \varepsilon \rangle + \frac{\sqrt{2} - \nu}{2} \langle d, \mu_{l} \rangle \right) \right) \right) \\ &\quad - \sum_{K_{1} + 1 \leq l \leq K} \exp \left(\frac{D}{\alpha^{2}} \left(\langle \mu_{i}, \varepsilon \rangle + \frac{\sqrt{2} - \nu}{2} \langle d, \mu_{l} \rangle \right) \right) \right) \\ &\leq \mathbb{P}_{\varepsilon} \left(\min_{\|d\| \leq 1} \left[\exp \left(\frac{D}{\alpha^{2}} \left(1 + \langle \mu_{k}, \varepsilon \rangle + \frac{\sqrt{2} - \nu}{2} \langle d, \mu_{k} \rangle \right) \right) \right) \\ &\quad - \sum_{K_{1} + 1 \leq l \leq K} \exp \left(\frac{D}{\alpha^{2}} \left(\langle \mu_{i}, \varepsilon \rangle + \frac{\sqrt{2} - \nu}{2} \langle d, \mu_{k} \rangle \right) \right) \right) \\ &\leq \mathbb{P}_{\varepsilon} \left(\min_{\|d\| \leq 1} \left[\exp \left(\frac{D}{\alpha^{2}} \left(1 + \langle \mu_{k}, \varepsilon \rangle + \frac{\sqrt{2} - \nu}{2} \langle d, \mu_{k} \rangle \right) \right) \right) \\ &\quad - \sum_{K_{1} + 1 \leq l \leq K} \exp \left(\frac{D}{\alpha^{2}} \left(\langle \mu_{i}, \varepsilon \rangle + \frac{\sqrt{2} - \nu}{2} \langle d, \mu_{k} \rangle \right) \right) \\ &\leq \mathbb{P}_{\varepsilon} \left(\min_{\|d\| \leq 1} \left[\exp \left(\frac{D}{\alpha^{2}} \left(1 + \langle \mu_{k}, \varepsilon \rangle + \frac{\sqrt{2} - \nu}{2} \langle d, \mu_{k} \rangle \right) \right) \right) \\ &\quad - \sum_{K_{1} + 1 \leq l \leq K} \exp \left(\frac{D}{\alpha^{2}} \left(|\langle \mu_{i}, \varepsilon \rangle| + \frac{\sqrt{2} - \nu}{2} |\langle d, \mu_{i} \rangle| \right) \right) \right) \right] \leq 0 \\ &\leq \mathbb{P}_{\varepsilon} \left(\min_{\|d\| \leq 1} \left[\exp \left(\frac{D}{\alpha^{2}} \left(1 + \langle \mu_{k}, \varepsilon \rangle + \frac{\sqrt{2} - \nu}{2} \langle d, \mu_{k} \rangle \right) \right) \\ &\quad - \sum_{K_{2} \exp \left(\frac{D}{\alpha^{2}} \left(1 + \langle \mu_{k}, \varepsilon \rangle + \frac{\sqrt{2} - \nu}{2} \langle d, \mu_{k} \rangle \right) \right) \\ &\leq \mathbb{P}_{\varepsilon} \left(\min_{\|d\| \leq 1} \left[1 + \langle \mu_{k}, \varepsilon \rangle + \frac{\sqrt{2} - \nu}{2} \langle d, \mu_{k} \rangle \right] \\ &\quad - \frac{\alpha^{2}}{D} \log K_{2} - \lim_{K_{1} + 1 \leq l \leq K} |\langle \mu_{i}, \varepsilon \rangle| + \frac{\sqrt{2} - \nu}{2} \lim_{K_{1} + 1 \leq l \leq K} |\langle d, \mu_{l} \rangle| \right) \right) \right] \leq 0 \\ &\leq \mathbb{P}_{\varepsilon} \left(\min_{\|d\| \leq 1} \left[1 + \frac{\sqrt{2} - \nu}{2} \langle d, \mu_{k} \rangle - \frac{\sqrt{2} - \nu}{2} \lim_{K_{1} + 1 \leq l \leq K} |\langle d, \mu_{l} \rangle| \right) \\ &\quad - \frac{\alpha^{2}}{D} \log K_{2} - |\langle \mu_{k}, \varepsilon \rangle| - \lim_{K_{1} + 1 \leq l \leq K} |\langle \mu_{i}, \varepsilon \rangle| \leq 0 \right), \quad (B.20) \\ \\ &\qquad \\ \lim_{\|d\| \leq 1} \left[1 + \frac{\sqrt{2} - \nu}{2} \sqrt{2} \sqrt{\left| \langle d, \mu_{k} \rangle|^{2} - \frac{\sqrt{2}}{2} \lim_{K_{1} + 1 \leq l \leq K} |\langle d, \mu_{l} \rangle|^{2}} \right] \\ &\geq \min_{\|d\| \leq 1} \left[1 - \frac{\sqrt{2} - \nu}{2} \sqrt{2} \sqrt{\left| \langle d, \mu_{k} \rangle|^{2} - \frac{\nu}{2} \lim_{K_{1} + 1 \leq l \leq K} |\langle \mu_{i}, \varepsilon \rangle|^{2} - \frac{\omega}{\sqrt{2}} \lim_{K_{1} + 1 \leq l \leq K} |\langle \mu_{i}, \varepsilon \rangle| \leq 0 \right) \\ \\ &\geq \min_{\|d\| \leq 1} \left[1 - \frac{\sqrt{2} - \nu}{2} \sqrt{2} \sqrt{\left| \langle d, \mu_{k} \rangle|^{2} - \frac{\omega^{2}}{2} \lim_{K_{1}$$

1134 1135	$\leq \mathbb{P}_{oldsymbol{arepsilon}}\left(rac{ u}{2\sqrt{2}}-2\max_{1\leq l\leq K} \langleoldsymbol{\mu}_l,oldsymbol{arepsilon} angle \leq 0 ight)$	
1136	$(\mu \nu \nu) = (\mu \nu^2)$	
1137	$\leq K \mathbb{P}_{\boldsymbol{\varepsilon}}\left(\left \langle \boldsymbol{\mu}_1, \boldsymbol{\varepsilon} angle ight \geq rac{1}{4\sqrt{2}} ight) \leq 2K \exp\left(-rac{1}{64lpha^2} ight) .$	(B.21)
1138		_
1139	The proof for the case $k \ge K_1 + 1$ is almost identical.	
1140		
11/10		
1143		
1144		
1145		
1146		
1147		
1148		
1149		
1150		
1151		
1152		
1153		
1154		
1155		
1156		
1157		
1158		
1159		
1160		
1161		
1162		
1103		
1165		
1166		
1167		
1168		
1169		
1170		
1171		
1172		
1173		
1174		
1175		
1176		
1177		
1178		
1179		
1180		
1100		
1102		
1103		
1185		
1186		
1187		

PRELU CONVERGES TO OPTIMAL ℓ_2 -ROBUST CLASSIFIER, PART ONE: С **CONVERGENCE IMPLIES ROBUSTNESS**

We prove Proposition 1 here.

Proposition 1 (Restated). Given a classifier f that satisfies $f(\gamma x) = \gamma f(x), \forall x \in \mathbb{R}^D, \forall \gamma > 0$ and $dist(f, F^{(p)}) = inf_{c>0} \sup_{x \in \mathbb{S}^{D-1}} |cf(x) - F^{(p)}(x)| \le \nu \text{ for some } p > 2 \text{ and } 0 < \nu \le \left(\frac{\sqrt{2}}{8}\right)^p.$ Then for a sample $(\boldsymbol{x}, y) \sim \mathcal{D}_{X,Y}$, we have

$$\mathbb{P}_{(\boldsymbol{x},\boldsymbol{y})\sim\mathcal{D}_{\boldsymbol{X},\boldsymbol{Y}}}\left(\min_{\|\boldsymbol{d}\|\leq 1}\left[f\left(\boldsymbol{x}+\frac{\sqrt{2}-8\nu^{\frac{1}{p}}}{2}\boldsymbol{d}\right)\boldsymbol{y}\right]>0\right)\geq 1-2K\exp\left(-\frac{D\nu^{\frac{2}{p}}}{2K^{2}\alpha^{2}}\right)-4\exp\left(-\frac{3}{8\alpha^{2}}\right).$$
(C.1)

Proof. First of all, since $f(\gamma x) = \gamma f(x), \forall x \in \mathbb{R}^D, \forall \gamma > 0$ and the same holds for $F^{(p)}(\cdot)$, we suppose the infimum is attained at $c^* \ge 0$, then

$$\sup_{\boldsymbol{x}\in\mathbb{R}^{D}}|c^{*}f(\boldsymbol{x})-F^{(p)}(\boldsymbol{x})|=\sup_{\boldsymbol{x}\in\mathbb{R}^{D}}\left|c^{*}f\left(\frac{\boldsymbol{x}}{\|\boldsymbol{x}\|}\right)-F^{(p)}\left(\frac{\boldsymbol{x}}{\|\boldsymbol{x}\|}\right)\right|\|\boldsymbol{x}\|\leq\|\boldsymbol{x}\|\nu\,,\qquad(C.2)$$

where the last inequality uses $dist(f, F^{(p)}) \leq \nu$. With (C.2), we have

$$\begin{split} & \mathbb{P}\left(\min_{\|d\|\leq 1} \left[f\left(x + \frac{\sqrt{2} - 8\nu^{\frac{1}{p}}}{2}d\right)y\right] \leq 0 \right) \\ & \mathbb{P}\left(\min_{\|d\|\leq 1} \left[c^{*}f\left(x + \frac{\sqrt{2} - 8\nu^{\frac{1}{p}}}{2}d\right)y\right] \leq 0 \right) \\ & \mathbb{P}\left(\min_{\|d\|\leq 1} \left[c^{*}f\left(x + \frac{\sqrt{2} - 8\nu^{\frac{1}{p}}}{2}d\right)y - F^{(p)}\left(x + \frac{\sqrt{2} - 8\nu^{\frac{1}{p}}}{2}d\right)y + F^{(p)}\left(x + \frac{\sqrt{2} - 8\nu^{\frac{1}{p}}}{2}d\right)y\right] \leq 0 \right) \\ & \mathbb{P}\left(\min_{\|d\|\leq 1} \left[c^{*}f\left(x + \frac{\sqrt{2} - 8\nu^{\frac{1}{p}}}{2}d\right)y - F^{(p)}\left(x + \frac{\sqrt{2} - 8\nu^{\frac{1}{p}}}{2}d\right)y - F^{(p)}\left(x + \frac{\sqrt{2} - 8\nu^{\frac{1}{p}}}{2}d\right)y - F^{(p)}\left(x + \frac{\sqrt{2} - 8\nu^{\frac{1}{p}}}{2}d\right)y\right] \leq 0 \right) \\ & \mathbb{P}\left(\min_{\|d\|\leq 1} \left[F^{(p)}\left(x + \frac{\sqrt{2} - 8\nu^{\frac{1}{p}}}{2}d\right)y\right] - \max_{\|d\|\leq 1} \left| c^{*}f\left(x)y - F^{(p)}\left(x\right)y\right| \leq 0, \ \|x\|^{2} \leq \frac{17}{2}\right) + \mathbb{P}\left(\|x\|^{2} > \frac{17}{2}\right) \\ & \leq \mathbb{P}\left(\min_{\|d\|\leq 1} \left[F^{(p)}\left(x + \frac{\sqrt{2} - 8\nu^{\frac{1}{p}}}{2}d\right)y\right] - \max_{\|x\|^{2}\leq 9} \left| c^{*}f\left(x\right)y - F^{(p)}\left(x\right)y\right| \leq 0, \ \|x\|^{2} \leq \frac{17}{2}\right) + \mathbb{P}\left(\|x\|^{2} > \frac{17}{2}\right) \\ & \leq \mathbb{P}\left(\min_{\|d\|\leq 1} \left[F^{(p)}\left(x + \frac{\sqrt{2} - 8\nu^{\frac{1}{p}}}{2}d\right)y\right] - \max_{\|x\|^{2}\leq 9} \left| c^{*}f\left(x\right)y - F^{(p)}\left(x\right)y\right| \leq 0 \right) + \mathbb{P}\left(\|x\|^{2} > \frac{17}{2}\right) \\ & \leq \mathbb{P}\left(\min_{\|d\|\leq 1} \left[F^{(p)}\left(x + \frac{\sqrt{2} - 8\nu^{\frac{1}{p}}}{2}d\right)y\right] \leq 3\nu \right) + \mathbb{P}\left(\|x\|^{2} > \frac{17}{2}\right) . \\ & \text{The second term } \mathbb{P}\left(\|x\|^{2} \geq \frac{17}{2}\right) \text{ is easy to bound our focus is to show} \\ \end{array}$$

The second term $\mathbb{P}\left(\|\boldsymbol{x}\|^2 > \frac{17}{2}\right)$ is easy to bound, our focus is to show

$$\mathbb{P}\left(\min_{\|\boldsymbol{d}\|\leq 1}\left[F^{(p)}\left(\boldsymbol{x}+\frac{\sqrt{2}-8\nu^{\frac{1}{p}}}{2}\boldsymbol{d}\right)\boldsymbol{y}\right]\leq 3\nu\right)\geq 2(K+1)\exp\left(-\frac{CD\nu^{2}}{K^{2}\alpha^{2}}\right),\qquad(C.3)$$

which resembles the result in Min & Vidal (2024, Theorem 1), but one can not directly obtain (C.3) from this existing result. Nonetheless, we can partially follow Min & Vidal (2024, Theorem 1)'s proof and obtain (C.3) (with non-trivial new derivations), as shown below:

Since

1237
1238
1239
1240
$$\mathbb{P}\left(\min_{\|\boldsymbol{d}\|\leq 1} \left[F^{(p)}\left(\boldsymbol{x} + \frac{\sqrt{2} - 8\nu^{\frac{1}{p}}}{2}\boldsymbol{d}\right)\boldsymbol{y}\right] \leq 3\nu\right)$$

$$K_{\boldsymbol{d}}\left(\sum_{\boldsymbol{d$$

1240
1241
$$= \sum_{k=1}^{K} \mathbb{P}\left(\min_{\|\boldsymbol{d}\| \le 1} \left[F^{(p)}\left(\boldsymbol{x} + \frac{\sqrt{2} - 8\nu^{\frac{1}{p}}}{2}\boldsymbol{d} \right) \boldsymbol{y} \right] \le 3\nu \left| \boldsymbol{z} = \boldsymbol{k} \right) \mathbb{P}\left(\boldsymbol{z} = \boldsymbol{k}\right)$$

1242
11 suffices to show that
$$\forall 1 \leq k \leq K$$

1244
1245
1246
1247
1248
1249
1249
1249
1249
1240
1240
1240
1240
1240
1240
1241
1240
1240
1240
1240
1240
1240
1240
1240
1251
1252
1253
1254
1255
1255
1255
1255
1255
1256
1257
1256
1257
1256
1257
1256
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1258
1257
1257
1258
1257
1257
1257
1257
1257
1257
1257
1257
1257
1257
1257
1257
1257
1257
1257
1257
1257
1257
1257
1257
1257
1257
1257
1257
1257
1257
1257
1257
1257
1257
1

$$\mathcal{E} := \left\{ 1 + \langle \boldsymbol{\mu}_k, \boldsymbol{\varepsilon} \rangle + \frac{\sqrt{2} - 8\nu^{\frac{1}{p}}}{2} \langle \boldsymbol{d}, \boldsymbol{\mu}_k \rangle \ge 0, \forall \boldsymbol{d} \in \mathbb{S}^{D-1} \right\},$$
(C.6)

1274 Then, by Min & Vidal (2024, Lemma 2), 1275

$$\begin{aligned} & \text{(C.5)} = \mathbb{P}_{\varepsilon} \left(\min_{\|\boldsymbol{d}\| \leq 1} \left[\sigma^{p} \left(1 + \langle \boldsymbol{\mu}_{k}, \varepsilon \rangle + \frac{\sqrt{2} - 8\nu^{\frac{1}{p}}}{2} \langle \boldsymbol{d}, \boldsymbol{\mu}_{k} \rangle \right) \right] \\ & \text{1278} \\ & \text{1279} \\ & \text{1280} \\ & \text{1280} \\ & \text{1281} \\ & \text{1282} \\ & \text{1282} \\ & \text{1283} \\ & \text{1284} \\ & \text{1284} \\ & \text{1284} \\ & \text{1284} \\ & \text{1285} \\ & \text{1286} \\ & \text{1286} \\ & \text{1286} \\ & \text{1287} \\ & \text{1287} \\ & \text{1287} \\ & \text{1288} \\ & \text{1288} \\ & \text{1286} \\ & \text{1287} \\ & \text{1288} \\ & \text{1289} \end{aligned}$$

1289
1290 Since under event
$$\mathcal{E}$$
, we have $\sigma^p \left(1 + \langle \boldsymbol{\mu}_k, \boldsymbol{\varepsilon} \rangle + \frac{\sqrt{2} - 8\nu^{\frac{1}{p}}}{2} \langle \boldsymbol{d}, \boldsymbol{\mu}_k \rangle \right) =$
1291
1292 $\left(1 + \langle \boldsymbol{\mu}_k, \boldsymbol{\varepsilon} \rangle + \frac{\sqrt{2} - 8\nu^{\frac{1}{p}}}{2} \langle \boldsymbol{d}, \boldsymbol{\mu}_k \rangle \right)^p$, we can proceed with
1293 (C.7) = $\mathbb{P}_{\boldsymbol{\varepsilon}} \left(\min_{\|\boldsymbol{d}\| \le 1} \left[\left(1 + \langle \boldsymbol{\mu}_k, \boldsymbol{\varepsilon} \rangle + \frac{\sqrt{2} - 8\nu^{\frac{1}{p}}}{2} \langle \boldsymbol{d}, \boldsymbol{\mu}_k \rangle \right)^p \right]$

$$\begin{aligned} & -\sum_{K_{1}+1\leq l\leq K}\sigma^{p}\left(\left(\mu,\varepsilon\right)+\frac{\sqrt{2}-8\nu^{\frac{1}{2}}}{2}\left\langle d,\mu_{l}\right\rangle\right)\right) \leq 3\nu, \mathcal{E}\right)+\mathbb{P}\left(\mathcal{E}^{c}\right) \\ &\leq \mathbb{P}_{\varepsilon}\left(\min_{\|d\|\leq 1}\left[\left(1+\left\langle\mu,\epsilon\right\rangle+\frac{\sqrt{2}-8\nu^{\frac{1}{2}}}{2}\left\langle d,\mu_{k}\right\rangle\right)^{p}\right. \\ & -\sum_{K_{1}+1\leq l\leq K}\left(\left\langle\left(\mu,\epsilon\right\rangle+\frac{\sqrt{2}-8\nu^{\frac{1}{2}}}{2}\right|\left\langle d,\mu_{l}\right\rangle\right)\right)^{p}-3\nu\right] < 0, \mathcal{E}\right)+\mathbb{P}\left(\mathcal{E}^{c}\right) \\ &\leq \mathbb{P}_{\varepsilon}\left(\min_{\|d\|\leq 1}\left[1+\left\langle\mu,\epsilon\right\rangle+\frac{\sqrt{2}-8\nu^{\frac{1}{2}}}{2}\left\langle d,\mu_{k}\right\rangle\right. \\ & -\left(\sum_{K_{1}+1\leq l\leq K}\left(\left|\left\langle\mu,\epsilon\right\rangle\right|+\frac{\sqrt{2}-8\nu^{\frac{1}{2}}}{2}\right|\left\langle d,\mu_{l}\right\rangle\right)\right)^{p}+3\nu\right)^{1/p}\right] < 0, \mathcal{E}\right)+\mathbb{P}\left(\mathcal{E}^{c}\right) \\ &\leq \mathbb{P}_{\varepsilon}\left(\min_{\|d\|\leq 1}\left[1+\frac{\sqrt{2}-8\nu^{\frac{1}{2}}}{2}\left\langle d,\mu_{k}\right\rangle-\frac{\sqrt{2}-8\nu^{\frac{1}{2}}}{2}\left\langle d,\mu_{l}\right\rangle\right]^{p}+3\nu\right)^{1/p}\right] < 0, \mathcal{E}\right)+\mathbb{P}\left(\mathcal{E}^{c}\right) \\ & -\left(\sum_{K_{1}+1\leq l\leq K}\left(\left|\left\langle\mu,\epsilon\right\rangle\right|+\frac{\sqrt{2}-8\nu^{\frac{1}{2}}}{2}\left\langle d,\mu_{k}\right\rangle\right)^{1/p}-\left|\left\langle\mu,\epsilon\right\rangle\right| < 0, \mathcal{E}\right)+\mathbb{P}\left(\mathcal{E}^{c}\right) \\ &= \mathbb{P}_{\varepsilon}\left(\min_{\|d\|\leq 1}\left[1+\frac{\sqrt{2}-8\nu^{\frac{1}{2}}}{2}\left\langle d,\mu_{k}\right\rangle-\frac{\sqrt{2}-8\nu^{\frac{1}{2}}}{2}\left\langle d,\mu_{l}\right\rangle\right]^{1/p}-\left|\left\langle\mu,\epsilon\right\rangle\right| < 0, \mathcal{E}\right)+\mathbb{P}\left(\mathcal{E}^{c}\right) \\ &= \mathbb{P}_{\varepsilon}\left(\max_{\|d\|\leq 1}\left[1+\frac{\sqrt{2}-8\nu^{\frac{1}{2}}}{2}\left\langle d,\mu_{k}\right\rangle\right] -\frac{\sqrt{2}-8\nu^{\frac{1}{2}}}{2}\left\langle d,\mu_{l}\right\rangle\right)^{1/p}-\left|\left\langle d,\nu\right\rangle\right) \\ &= \mathbb{P}_{\varepsilon}\left(\sum_{K_{1}+1\leq l\leq K}\left(\left|\left\langle\mu,\mu,\varepsilon\right\rangle\right|+\frac{\sqrt{2}-8\nu^{\frac{1}{2}}}{2}\left\langle d,\mu_{k}\right\rangle\right)\right)^{1/p}-\left|\left\langle\mu,\mu,\varepsilon\right\rangle\right| < 0, \mathcal{E}\right) +\mathbb{P}\left(\mathcal{E}^{c}\right) \\ &= \mathbb{P}_{\varepsilon}\left(\sum_{K_{1}+1\leq l\leq K}\left|\left\langle\mu,\mu,\varepsilon\right\rangle\right|+\left\langle\mu,\nu\right\rangle\right| < \frac{\sqrt{2}}{2}\left\langle d,\mu_{k}\right\rangle\right) \\ &= \mathbb{P}_{\varepsilon}\left(\sum_{K_{1}+1\leq l\leq K}\left|\left\langle\mu,\mu,\varepsilon\right\rangle\right|+\left\langle\mu,\nu\right\rangle\right| < 2\sqrt{2}\nu^{\frac{1}{2}}\right) +\mathbb{P}\left(\mathcal{E}^{c}\right) \\ &\geq \mathbb{P}_{\varepsilon}\left(\sum_{K_{1}+1\leq l\leq K}\left|\left\langle\mu,\mu,\varepsilon\right\rangle\right|+\left|\left\langle\mu,\nu\right\rangle\right| > \sqrt{2}\nu^{\frac{1}{2}}\right) +\mathbb{P}\left(\mathcal{E}^{c}\right) \\ &\geq \mathbb{P}_{\varepsilon}\left(\sum_{K_{1}+1\leq l\leq K}\left|\left\langle\mu,\mu,\varepsilon\right\rangle\right|+\left\langle\mu,\nu\right\rangle\right| > \sqrt{2}\nu^{\frac{1}{2}}\right) +\mathbb{P}\left(\mathcal{E}^{c}\right) \\ &\geq \mathbb{P}_{\varepsilon}\left(\sum_{K_{1}+1\leq l\leq K}\left|\left\langle\mu,\mu,\varepsilon\right\rangle| > \frac{\sqrt{2}\nu^{\frac{1}{2}}}{K}\right) +\mathbb{P}\left(\mathcal{E}^{c}\right) \\ &\geq \mathbb{P}_{\varepsilon}\left(\sum_{K_{1}+1\leq l\leq K}\left|\left\langle\mu,\mu,\varepsilon\right\rangle\right| > \frac{\sqrt{2}\nu^{\frac{1}{2}}}{K}\right) +\mathbb{P}\left(\mathcal{E}^{c}\right) \\ &\geq \mathbb{P}_{\varepsilon}\left(1+|\mu,\varepsilon\rangle\right) > \frac{\sqrt{2}\nu^{\frac{1}{2}}}{K}\right) +\mathbb{P}\left(\mathcal{E}^{c}\right) \\ &\geq \mathbb{P}\left(\left|\left\langle\mu,\nu\right\rangle\right| > \frac{\sqrt{2}\nu^{\frac{1$$

 $\mathbb{P}\left(\|\boldsymbol{x}\|^2 > \frac{17}{2}\right) \leq \mathbb{P}\left(\|\boldsymbol{\varepsilon}\| \geq \sqrt{\frac{17}{2}} - 1\right) \leq 4\exp\left(-\left(\sqrt{\frac{17}{2}} - 1\right)^2 \frac{1}{8\alpha^2}\right) \leq 4\exp\left(-\frac{3}{8\alpha^2}\right),$

The proof is finished, notice that the bad event $||x||^2 > \frac{17}{2}$ is chosen arbitrarily, so one can derive more general results by letting the results depend on the choice of a bad event. But for our purpose, we do not need it.

1404
1405
1406DPRELU CONVERGES TO OPTIMAL ℓ_2 -ROBUST CLASSIFIER, PART TWO:
BASIC RESULTS ON NEURON DYNAMICS AND GOOD EVENTS

We also let $\mathcal{I}_k := \{i : (k-1)N + 1 \le i \le kN\}$, the index set of data sampled from k-th cluster.

1417 D.1 RESULTS ON NEURON DYNAMICS

Neuron dynamics: Under GF, we have

$$\frac{d}{dt}\boldsymbol{w}_{j} = -\frac{1}{N}\sum_{i=1}^{KN} \nabla_{\hat{y}}\ell_{i} v_{j} \left(\frac{p[\sigma(\langle\boldsymbol{x}_{i},\boldsymbol{w}_{j}\rangle)]^{p-1}}{\|\boldsymbol{w}_{j}\|^{p-1}}\boldsymbol{x}_{i} - (p-1)\frac{[\sigma(\langle\boldsymbol{x}_{i},\boldsymbol{w}_{j}\rangle)]^{p}}{\|\boldsymbol{w}_{j}\|^{p+1}}\boldsymbol{w}_{j}\right)$$

$$= -\frac{1}{N}\sum_{i:\langle\boldsymbol{x}_{i},\boldsymbol{w}_{j}\rangle>0} \nabla_{\hat{y}}\ell_{i} v_{j} \left(\frac{p[\langle\boldsymbol{x}_{i},\boldsymbol{w}_{j}\rangle]^{p-1}}{\|\boldsymbol{w}_{j}\|^{p-1}}\boldsymbol{x}_{i} - (p-1)\frac{[\langle\boldsymbol{x}_{k},\boldsymbol{w}_{j}\rangle]^{p}}{\|\boldsymbol{w}_{j}\|^{p+1}}\boldsymbol{w}_{j}\right)$$

and similarly,

$$\frac{d}{dt}v_j = -\frac{1}{N} \sum_{i: \langle \boldsymbol{x}_i, \boldsymbol{w}_j \rangle > 0} \nabla_{\hat{y}} \ell_i \frac{[\langle \boldsymbol{x}_i, \boldsymbol{w}_j \rangle]^p}{\|\boldsymbol{w}_j\|^{p-1}}$$

Balancedness: We compute

$$\begin{split} \frac{d}{dt}(\boldsymbol{w}_{j}^{\top}\boldsymbol{w}_{j}) &= 2\left\langle \frac{d}{dt}\boldsymbol{w}_{j}, \boldsymbol{w}_{j} \right\rangle \\ &= -\frac{1}{N}\sum_{i:\langle \boldsymbol{x}_{i}, \boldsymbol{w}_{j} \rangle > 0} \nabla_{\hat{y}}\ell_{i} v_{j} \left(\frac{p[\langle \boldsymbol{x}_{i}, \boldsymbol{w}_{j} \rangle]^{p}}{\|\boldsymbol{w}_{j}\|^{p-1}} - (p-1)\frac{[\langle \boldsymbol{x}_{k}, \boldsymbol{w}_{j} \rangle]^{p}}{\|\boldsymbol{w}_{j}\|^{p-1}} \right) \\ &= -\frac{2}{N}\sum_{i:\langle \boldsymbol{x}_{i}, \boldsymbol{w}_{j} \rangle > 0} \nabla_{\hat{y}}\ell_{i} v_{j} \frac{[\langle \boldsymbol{x}_{i}, \boldsymbol{w}_{j} \rangle]^{p}}{\|\boldsymbol{w}_{j}\|^{p-1}} \,, \end{split}$$

1441 and

$$\frac{d}{dt}v_j^2 = 2 v_j \frac{d}{dt}v_j = -\frac{2}{N} \sum_{i: \langle \boldsymbol{x}_i, \boldsymbol{w}_j \rangle > 0} \nabla_{\hat{y}} \ell_i v_j \frac{[\langle \boldsymbol{x}_i, \boldsymbol{w}_j \rangle]^p}{\|\boldsymbol{w}_j\|^{p-1}}$$

1445 Therefore, we have

$$\frac{d}{dt}(\boldsymbol{w}_j^{\top}\boldsymbol{w}_j - v_j^2) \equiv 0, \qquad (D.1)$$

thus $\boldsymbol{w}_{j}^{\top}(t)\boldsymbol{w}_{j}(t) - v_{j}^{2}(t) = \boldsymbol{w}_{j}^{\top}(0)\boldsymbol{w}_{j}(0) - v_{j}^{2}(0), \forall t$, since we have a balanced initialization such that $\boldsymbol{w}_{j}^{\top}(0)\boldsymbol{w}_{j}(0) - v_{j}^{2}(0), \forall j$. Such balancedness holds for all time t. Using this balancedness $v_{j}^{2} \equiv \|\boldsymbol{w}_{j}\|^{2}, \forall j \in [h]$, we can write

$$\frac{1452}{1453} \qquad \frac{d}{dt} \boldsymbol{w}_{j} = -\frac{\operatorname{sign}(v_{j}(0))}{N} \sum_{i:\langle \boldsymbol{x}_{i}, \boldsymbol{w}_{j} \rangle > 0} \nabla_{\hat{y}} \ell_{i} \|\boldsymbol{w}_{j}\| \left(p\left(\left\langle \boldsymbol{x}_{i}, \frac{\boldsymbol{w}_{j}}{\|\boldsymbol{w}_{j}\|} \right\rangle \right)^{p-1} \boldsymbol{x}_{i} - (p-1)\left(\left\langle \left\langle \boldsymbol{x}_{i}, \frac{\boldsymbol{w}_{j}}{\|\boldsymbol{w}_{j}\|} \right\rangle \right)^{p} \frac{\boldsymbol{w}_{j}}{\|\boldsymbol{w}_{j}\|} \right)$$

$$(D.2)$$

where we use that $sign(v_j(t)) = sign(v_j(0))$, which is another consequence of balancedness Boursier et al. (2022); Min et al. (2024). We will study the dynamics of w_j from now on, and one can write the time derivatives of the norm and direction of these neurons:

Neuron angular dynamics:

Neuron norm dynamics:

 $= 2\left\langle \boldsymbol{w}_{j}, \frac{d}{dt}\boldsymbol{w}_{j} \right\rangle$

 $\frac{d}{dt} \| \boldsymbol{w}_j \|^2$

$$\frac{d}{dt} \frac{\boldsymbol{w}_{j}}{\|\boldsymbol{w}_{j}\|} = \left(I - \frac{\boldsymbol{w}_{j}\boldsymbol{w}_{j}^{\top}}{\|\boldsymbol{w}_{j}\|^{2}}\right) \frac{1}{\|\boldsymbol{w}_{j}\|} \frac{d}{dt} \boldsymbol{w}_{j} = -\frac{\operatorname{sign}(v_{j}(0))}{N} \sum_{i:\langle \boldsymbol{x}_{i}, \boldsymbol{w}_{j} \rangle > 0} \nabla_{\hat{y}} \ell_{i} \left(I - \frac{\boldsymbol{w}_{j}\boldsymbol{w}_{j}^{\top}}{\|\boldsymbol{w}_{j}\|^{2}}\right) \left(p\left(\left\langle\boldsymbol{x}_{i}, \frac{\boldsymbol{w}_{j}}{\|\boldsymbol{w}_{j}\|}\right\rangle\right)^{p-1} \boldsymbol{x}_{i} - (p-1)\left(\left\langle\boldsymbol{x}_{i}, \frac{\boldsymbol{w}_{j}}{\|\boldsymbol{w}_{j}\|}\right\rangle\right)^{p} \frac{\boldsymbol{w}_{j}}{\|\boldsymbol{w}_{j}\|}\right) = -\frac{\operatorname{sign}(v_{j}(0))}{N} \sum_{i:\langle \boldsymbol{x}_{i}, \boldsymbol{w}_{j} \rangle > 0} \nabla_{\hat{y}} \ell_{i} p\left(\left\langle\boldsymbol{x}_{i}, \frac{\boldsymbol{w}_{j}}{\|\boldsymbol{w}_{j}\|}\right\rangle\right)^{p-1} \left(\boldsymbol{x}_{i} - \left\langle\boldsymbol{x}_{i}, \frac{\boldsymbol{w}_{j}}{\|\boldsymbol{w}_{j}\|}\right\rangle \frac{\boldsymbol{w}_{j}}{\|\boldsymbol{w}_{j}\|}\right). \quad (D.4)$$

 $= -2 \frac{\operatorname{sign}(v_j(0))}{N} \sum_{i: \langle \boldsymbol{x}_i, \boldsymbol{w}_j \rangle > 0} \nabla_{\hat{y}} \ell_i \| \boldsymbol{w}_j \| \left(p\left(\left\langle \boldsymbol{x}_i, \frac{\boldsymbol{w}_j}{\|\boldsymbol{w}_j\|} \right\rangle \right)^{p-1} \langle \boldsymbol{w}_j, \boldsymbol{x}_i \rangle - (p-1) \left(\left\langle \left\langle \boldsymbol{x}_i, \frac{\boldsymbol{w}_j}{\|\boldsymbol{w}_j\|} \right\rangle \right)^p \| \boldsymbol{w}_j \| \right) \right)^{p-1} \langle \boldsymbol{w}_j, \boldsymbol{x}_j \rangle = 0$

(D.3)

 $= -2 \frac{\operatorname{sign}(v_j(0))}{N} \sum_{i: \langle \boldsymbol{x}_i, \boldsymbol{w}_i \rangle > 0} \nabla_{\hat{y}} \ell_k \|\boldsymbol{w}_j\| \left(p\left(\left\langle \boldsymbol{x}_i, \frac{\boldsymbol{w}_j}{\|\boldsymbol{w}_j\|} \right\rangle \right)^p \|\boldsymbol{w}_j\| - (p-1) \left(\left\langle \left\langle \boldsymbol{x}_i, \frac{\boldsymbol{w}_j}{\|\boldsymbol{w}_j\|} \right\rangle \right)^p \|\boldsymbol{w}_j\| \right)$

Finally, from the directional dynamics $\frac{d}{dt} \frac{\boldsymbol{w}_j}{\|\boldsymbol{w}_j\|}$, we obtain

 $= -2 \frac{\operatorname{sign}(v_j(0))}{N} \left(\sum_{i: (\boldsymbol{x}_i, \boldsymbol{w}_i) \geq 0} \nabla_{\hat{y}} \ell_i \left(\left\langle \boldsymbol{x}_i, \frac{\boldsymbol{w}_j}{\|\boldsymbol{w}_j\|} \right\rangle \right)^p \right) \|\boldsymbol{w}_j\|^2$

$$\frac{d}{dt}c_{kj} = \left\langle \boldsymbol{\mu}_{k}, \frac{d}{dt} \frac{\boldsymbol{w}_{j}}{\|\boldsymbol{w}_{j}\|} \right\rangle$$

$$= -\frac{\operatorname{sign}(v_{j}(0))}{N} \sum_{i:\langle \boldsymbol{x}_{i}, \boldsymbol{w}_{j} \rangle > 0} \nabla_{\hat{y}} \ell_{i} p\left(\left\langle \boldsymbol{x}_{i}, \frac{\boldsymbol{w}_{j}}{\|\boldsymbol{w}_{j}\|} \right\rangle\right)^{p-1} \left(\left\langle \boldsymbol{\mu}_{k}, \boldsymbol{x}_{i} \right\rangle - \left\langle \boldsymbol{x}_{i}, \frac{\boldsymbol{w}_{j}}{\|\boldsymbol{w}_{j}\|} \right\rangle c_{kj}\right),$$
(D.5)

and whenever $|c_{kj}| \neq 0$, we have

$$\frac{d}{dt} \log |c_{kj}| = \frac{1}{c_{kj}} \frac{d}{dt} c_{kj} = -\frac{\operatorname{sign}(v_j(0))}{N} \sum_{i:\langle \boldsymbol{x}_i, \boldsymbol{w}_j \rangle > 0} \nabla_{\hat{y}} \ell_i \, p\left(\left\langle \boldsymbol{x}_i, \frac{\boldsymbol{w}_j}{\|\boldsymbol{w}_j\|} \right\rangle\right)^{p-1} \left(\frac{\langle \boldsymbol{\mu}_k, \boldsymbol{x}_i \rangle}{c_{kj}} - \left\langle \boldsymbol{x}_i, \frac{\boldsymbol{w}_j}{\|\boldsymbol{w}_j\|} \right\rangle\right) \quad (D.6)$$

Our proof has the same structure as prior works Boursier et al. (2022); Min et al. (2024): We will
study neuron's angular dynamics (D.5) at the early phase (alignment phase) of the GF training, and
then study neuron's norm dynamics (D.3) at the later phase (convergence phase).

Lastly, in order to prove Lemma 7 and Proposition 2 in the next subsection, we need the following:

1511 We let $\{\mu_{K+1}, \dots, \mu_D\}$ be an orthonormal basis for the subspace that is orthogonal to $\operatorname{span}\{\mu_1, \dots, \mu_K\}$, and we can define $c_{kj} = \cos(\mu_k, w_j), k = K+1, \dots, D$. Since $\{\mu_1, \dots, \mu_D\}$

forms an orthonormal basis for the ambient space \mathbb{R}^D , we have

$$\sum_{k=1}^{D} c_{kj}^2 = \sum_{k=1}^{D} \left| \left\langle \boldsymbol{\mu}_k, \frac{\boldsymbol{w}_j}{\|\boldsymbol{w}_j\|} \right\rangle \right|^2 = 1.$$
 (D.7)

Moreover, we can write the same time-derivatives $\frac{d}{dt}c_{kj}$, $\frac{d}{dt}\log|c_{kj}|$ for $c_{kj} = \cos(\mu_k, w_j)$, k = $K + 1, \dots, D$ as in (D.5) and (D.6), respectively.

Lastly, the following inequality will be used frequently in our proof:

$$\sum_{l \neq k} c_{lj}^p \le \sum_{1 \le l \le D, l \neq k} |c_{lj}|^p \le \left(\sum_{1 \le l \le D, l \neq k} c_{lj}^2\right)^{\frac{p}{2}} = \left(1 - c_{kj}^2\right)^{\frac{p}{2}}$$
(D.8)

Note: The sum operation $\sum_{l \neq k}$ implicitly assumes $l \leq K$. We will explicitly indicate the range of lif it can take values between K + 1 and D.

D.2 GOOD EVENT

For a balanced dataset $\hat{\mathcal{D}} = \{ \boldsymbol{x}_i, y_i \}_{i=1}^{KN}$, notice that $\boldsymbol{x}_i = \boldsymbol{\mu}_{\lceil \frac{i}{N} \rceil} + \boldsymbol{\varepsilon}_i$ for some $\boldsymbol{\varepsilon}_i \in \mathcal{N}\left(0, \frac{\alpha^2}{D}\boldsymbol{I}\right)$. We define the following good event w.r.t. these ε_i s and show that they happen with high probability: **Lemma 4.** We define the event \mathcal{E}_{good} when the following happens:

1535
1536
$$I. \|\boldsymbol{\varepsilon}_i\| \leq \sqrt{8 \log \frac{16KN}{\delta}} \alpha, \ \forall 1 \leq i \leq KN;$$

1537

1537
1538
2.
$$|\langle \boldsymbol{\mu}_{k}, \boldsymbol{\varepsilon}_{i} \rangle| \leq \sqrt{2 \log \frac{8K^{2}N}{\delta}} \frac{\alpha}{\sqrt{D}}, \forall 1 \leq i \leq KN, 1 \leq k \leq K;$$

1539
1540
3. $\|\sum_{i \in \mathcal{I}_{k}} \boldsymbol{\varepsilon}_{i}\| \leq \sqrt{2 \log \frac{8K}{\delta}} \alpha \sqrt{N}, \forall 1 \leq k \leq K$
1541
4. $\sum_{i \in \mathcal{I}_{k}} \|\boldsymbol{\varepsilon}_{i}\|^{2} \leq 8 \log \frac{16K}{\delta} \alpha^{2}N, \forall 1 \leq k \leq K$
1542
4. $\sum_{i \in \mathcal{I}_{k}} \|\boldsymbol{\varepsilon}_{i}\|^{2} \leq 8 \log \frac{16K}{\delta} \alpha^{2}N, \forall 1 \leq k \leq K$
1544
1545
1. $\|\boldsymbol{\varepsilon}_{i}\| \leq C\sqrt{\log \frac{K^{2}N}{\delta}} \alpha, \forall 1 \leq i \leq KN;$
1546
1. $\|\boldsymbol{\varepsilon}_{i}\| \leq C\sqrt{\log \frac{K^{2}N}{\delta}} \alpha, \forall 1 \leq i \leq KN;$
1547
1. $\|\boldsymbol{\varepsilon}_{i}\| \leq C\sqrt{\log \frac{K^{2}N}{\delta}} \alpha, \forall 1 \leq i \leq KN, 1 \leq k \leq K;$
1550
1551
3. $\|\sum_{i \in \mathcal{I}_{k}} \boldsymbol{\varepsilon}_{i}\| \leq C\sqrt{\log \frac{K}{\delta}} \alpha \sqrt{N}, \forall 1 \leq k \leq K;$
1553
4. $\sum_{i \in \mathcal{I}_{k}} \|\boldsymbol{\varepsilon}_{i}\|^{2} \leq C \log \frac{K}{\delta} \alpha^{2}N, \forall 1 \leq k \leq K,$
1554
1555
1556
1557
Proof. We proof relevant probabilities one by one:

Proof. We proof relavent probabilities one by one:

1. By Lemma 3, we have

$$\mathbb{P}\left(\left\|\boldsymbol{\varepsilon}_{i}\right\| \geq t\right) \leq 4\exp\left(-\frac{t^{2}}{8\alpha^{2}}\right).$$
(D.9)

2. By Lemma 2, we have

$$\mathbb{P}\left(\left|\left\langle \boldsymbol{\mu}_{k}, \boldsymbol{\varepsilon}_{i}\right\rangle\right| \geq t\right) \leq 2\exp\left(-\frac{Dt^{2}}{2\alpha^{2}}\right).$$
(D.10)

1566 3. Apply Lemma 3 to the vector $\sum_{i \in \mathcal{I}_k} \varepsilon_i$, we have

$$\mathbb{P}\left(\left\|\sum_{i\in\mathcal{I}_{k}}\varepsilon_{i}\right\|\geq t\right)\leq4\exp\left(-\frac{t^{2}}{8N\alpha^{2}}\right).$$
(D.11)

Apply Lemma 3 to the vector that is the concatenation of all ε_i, i ∈ N_k and notice that its norm is equal to √∑_{i∈I_k} ||ε_i||², hence

$$\mathbb{P}\left(\sum_{i\in\mathcal{I}_k}\|\boldsymbol{\varepsilon}_i\|^2 \ge t^2\right) \le 4\exp\left(-\frac{t^2}{8N\alpha^2}\right).$$
(D.12)

Therefore,

$$\mathbb{P}\left(\|\boldsymbol{\varepsilon}_i\| \ge \sqrt{8\log\frac{16KN}{\delta}}\alpha\right) \le \frac{\delta}{4KN}, \qquad \forall 1 \le i \le KN,$$

$$\mathbb{P}\left(\left|\left\langle \boldsymbol{\mu}_{k},\boldsymbol{\varepsilon}_{i}\right\rangle\right| \geq \sqrt{2\log\frac{8K^{2}N}{\delta}}\frac{\alpha}{\sqrt{D}}\right) \leq \frac{\delta}{4K^{2}N}, \qquad \forall 1 \leq i \leq KN, 1 \leq k \leq K,$$

$$\mathbb{P}\left(\left\|\sum_{i\in\mathcal{I}_k}\boldsymbol{\varepsilon}_i\right\| \geq \sqrt{8\log\frac{16K}{\delta}}\alpha\sqrt{N}\right) \leq \frac{\delta}{4K}, \qquad \forall 1\leq k\leq K\,,$$

$$\mathbb{P}\left(\sum_{i\in\mathcal{I}_k}\|\boldsymbol{\varepsilon}_i\|^2 \ge 8\log\frac{16K}{\delta}\alpha^2 N\right) \le 4\exp\left(-\frac{t^2}{8N\alpha^2}\right) \le \frac{\delta}{4K}, \qquad \forall 1\le k\le K.$$

The union bound shows that $\mathbb{P}\left(\mathcal{E}_{good}\right) \leq 1 - \delta$.

1620EPRELU CONVERGES TO OPTIMAL ℓ_2 -ROBUST CLASSIFIER, PART THREE:1621ALIGNMENT PHASE

1623 1624 E.1 AUXILIARY LEMMAS

¹⁶²⁵ We need the following lemmas (proofs provided in Appendix G)

Lemma 5. Given an initialization shape that satisfies Assumption 2 with non-degeneracy gap $\Delta > 0$, then for $j \in \mathcal{N}_k$, we have

1629

$$c_{kj}(0) = \cos(\boldsymbol{\mu}_k, \boldsymbol{w}_j(0)) \ge \sqrt{\frac{1}{2} \left(\frac{1}{(1-\Delta)^2} - 1\right)} := \tilde{\Delta}_1, \tag{E.1}$$

1630 1631 1632

1633

1636

1637 1638 1639

1642 1643

1647 1648

1651 1652 1653

1656

1658 1659

1663 1664 1665

1670

1673

 $\frac{c_{lj}^{p-2}(0)}{c_{kj}^{p-2}(0)} \le (1 - \sqrt{2\Delta})^{p-2} := 1 - \tilde{\Delta}_2, \forall l \neq k \text{ with } y_l = y_k \text{ and } c_{lj}(0) > 0$ (E.2)

Lemma 6. Let p > 2. Condition on good event \mathcal{E}_{good} . Given some $1 \le k \le K$ and some $j \in \mathcal{N}_k$ and suppose the following is true at some point on the GF trajectory:

1.
$$c_{kj} \geq \tilde{\Delta}_1$$
;

2.
$$\frac{|c_{lj}|}{c_{kj}} \le (1 - \sqrt{2\Delta}), \forall l \neq k.$$

1640 *Then the following holds:* 1641

$$\frac{d}{dt}c_{kj} \ge pc_{kj}^{p-1}\tilde{\Delta}_2(1-c_{kj}) - C_1\log\frac{K}{\delta}\alpha^2 - C_2\max_k |f^{(p)}(\boldsymbol{\mu}_k;\boldsymbol{\theta}(t))|,$$

for some universal constant C_1, C_2 that depends on p. If one further assume $c_{kj} \ge \sqrt{\frac{4}{5}}$, then the lower bound can be improved as

$$\frac{d}{dt}c_{kj} \ge pc_{kj}^{p-1}\left(1 - \frac{1}{2^{p-2}}\right)(1 - c_{kj}) - C_1 \log \frac{K}{\delta}\alpha^2 - C_2 \max_i |f^{(p)}(\boldsymbol{x}_i; \boldsymbol{\theta})|$$

Lemma 7. Let p > 2. Condition on good event \mathcal{E}_{good} . Given an initialization shape that satisfies Assumption 2 with non-degeneracy gap $\Delta > 0$, define

$$t_{1a} := \inf\left\{t: \max_{i} |f^{(p)}(\boldsymbol{x}_{i}; \boldsymbol{\theta}(t)| > \min\left\{\frac{\tilde{\Delta}_{1}^{p-1}\tilde{\Delta}_{2}(1-\tilde{\Delta}_{1})}{2^{p+1}}, \frac{\tilde{\Delta}_{1}^{p-1}\tilde{\Delta}_{2}(1-\sqrt{2\Delta})}{2K2^{p+1}}\right\}\right\}.$$
(E.3)

1654 Then the following holds $\forall t \leq t_{1a}$: 1655

$$c_{kj}(t) \ge c_{kj}(0) \ge \tilde{\Delta}_1, \forall 1 \le k \le K, j \in \mathcal{N}_k ,$$
(E.4)

1657 and

$$\frac{|c_{lj}^{p-2}(t)|}{c_{kj}^{p-2}(t)} \le \frac{|c_{lj}^{p-2}(0)|}{c_{kj}^{p-2}(0)} \le 1 - \tilde{\Delta}_2 \text{ and } \forall l \neq k, j \in \mathcal{N}_k \text{.}$$
(E.5)

Lemma 8. Let p > 2. Condition on good event \mathcal{E}_{good} , then with any balanced initialization scale $\epsilon \leq \frac{1}{4\sqrt{h}W_{max}^2}$, the solution to gradient flow dynamics satisfies

$$\max_{k} |f^{(p)}(\boldsymbol{\mu}_{k};\boldsymbol{\theta}(t))| \leq 2\epsilon \sqrt{h} W_{\max}^{2}, \quad \forall t \leq \frac{1}{2^{p+2}K} \log\left(\frac{1}{2^{p-1}\sqrt{h\epsilon}}\right).$$
(E.6)

The following lemma will be used to upper-bound the time each neuron spends until reaching a neighborhood of some data μ_k .

Lemma 9. Let p > 2. Given some C > 0, if for some z(t), the following holds

$$\frac{d}{dt}z \ge Cz^{p-1}, \forall t \in [0,T], \ z(0) = z_0, \ z(T) = z_1,$$
(E.7)

for some $0 < z_0 \le z_1 < 1$. Then the travel time T for z(t) to go from z_0 to z_1 satisfies:

$$T \le \frac{1}{(p-2)Cz_0^{p-2}}.$$
(E.8)

Lemma 10. Let p > 2. Given some C > 0, if for some z(t), the following holds

$$\frac{d}{dt}z \ge C(1-z), \forall t \in [0,T], \ z(0) = z_0, \ z(T) = z_1,$$
(E.9)

for some $0 < z_0 \le z_1 < 1$. Then the travel time T for z(t) to go from z_0 to z_1 satisfies:

$$T \le \frac{1}{C} \log \frac{1}{1 - z_1}$$
 (E.10)

The following lemma will be used to lower-bound the time each neuron can stay around the neighborhood of some data μ_k .

1686 E.2 PROOF OF PROPOSITION 2

Proposition 2 (Restated). Given the same assumptions as in Theorem 3 and consider the same GF solution $\theta(t), t \ge 0$. There exist some $t_1 = \mathcal{O}\left(\log \frac{1}{\alpha}\right)$ and $t_2 = \mathcal{O}\left(\log \frac{1}{\epsilon}\right)$ such that $\forall k$ and $\forall j \in \mathcal{N}_k, \cos(\mu_k, w_j(t)) \ge 1 - \tilde{\mathcal{O}}(\alpha^2), \forall t \in [t_1, t_2]$.

1692 *Proof of Proposition 2.* Breakdown the proofs We let

$$t_1 := \inf \left\{ t : \min_k \min_{j \in \mathcal{N}_k} c_{kj}(t) \ge 1 - C \log \frac{K}{\delta} \alpha^2 \right\}.$$
 (E.11)

1696 We define

1676 1677

1680 1681 1682

1685

1687

1691

1693 1694 1695

$$\begin{aligned} & \mathbf{f}_{1697} \\ & \mathbf{f}_{0}_{1698} \\ & \mathbf{f}_{0} := \min \left\{ \frac{\tilde{\Delta}_{1}^{p-1} \tilde{\Delta}_{2} (1 - \tilde{\Delta}_{1})}{2^{p+2} \sqrt{h} W_{\max}^{2}}, \\ & \frac{\tilde{\Delta}_{1}^{p-1} \tilde{\Delta}_{2} (1 - \sqrt{2\Delta})}{2K 2^{p+2} \sqrt{h} W_{\max}^{2}}, \\ & \frac{p \tilde{\Delta}_{1}^{p-1} \tilde{\Delta}_{2} \alpha^{2}}{8 \sqrt{h} W_{\max}^{2}}, \\ & \frac{p \tilde{\Delta}_{1}^{p-1} \tilde{\Delta}_{2} \alpha^{2}}{8 \sqrt{h} W_{\max}^{2}}, \\ & \frac{1}{\sqrt{h}} \exp \left(-4K \left(\frac{20}{(p-2)p \tilde{\Delta}_{2} \tilde{\Delta}_{1}^{p-2}} + \frac{2}{p(2^{p-1}-2)} \log \frac{1}{C \log \frac{K}{\delta} \alpha^{2}} \right) \right) \right\}. \end{aligned}$$
(E.12)

Our goal is to show that if the initialization scale $\epsilon \leq \epsilon_0$ (Notice that our assumption $\epsilon = \Theta(\alpha^{8K})$ can satisfies this inequality), then

1712 1713 1. $\min_k \min_{j \in \mathcal{N}_k} c_{kj}(t)$ grows above $1 - C \log \frac{K}{\delta} \alpha^2$ before 1714 $\bar{t}_1 := \frac{20}{(p-2)p\tilde{\Delta}_2\tilde{\Delta}_1^{p-2}} + \frac{2}{p(2^{p-1}-2)} \log \frac{1}{C \log \frac{K}{\delta} \alpha^2};$

1716 2. Any $c_{kj}(t)$ staying above $1 - C \log \frac{K}{\delta} \alpha^2$ during $[t_1, t_2]$, where $t_2 := \frac{1}{2^{p+2}K} \log \left(\frac{1}{2^{p-1}\sqrt{h\epsilon}}\right)$; 1718

1719 The remaining proof is to show them one by one.

1720 1721 1722 1723 Upper bound on t_1 When $1 \le k \le K_1$, $j \in \mathcal{N}_k$ implies that $w_{j0} \in \mathcal{R}_k$ and $\operatorname{sign}(v_j) = 1$. We shall primarily focus on this case as the proof is nearly identical for $K_1 + 1 \le k \le K$. We prove it by contradiction.

1724 $\forall t \leq \bar{t}_1$, we have

1728 and

1732

1735

1736 1737

1739

1745 1746 1747

$$c_{kj}(t) \ge \tilde{\Delta}_1, \ \frac{c_{lj}^{p-2}}{c_{kj}^{p-2}} \le 1 - \tilde{\Delta}_2, \forall l \ne k, j \in \mathcal{N}_k.$$
(By (E.13) and Lemma 7) (E.14)

1733 Suppose $t_1 \ge \bar{t}_1$, then $\exists k, j \in \mathcal{N}_k$ such that $t_{1j}^{(k)} := \inf\{t : c_{kj}(t) \ge 1 - \frac{\alpha^2}{2}\} > \bar{t}_1$. However, for 1734 $0 \le t \le \bar{t}_1$, we have, by Lemma 6, for this particular k, j,

Whenever
$$c_{kj} \ge \tilde{\Delta}_1$$
,
 $\frac{d}{dt}c_{kj} \ge pc_{kj}^{p-1}\tilde{\Delta}_2(1-c_{kj}) - C_1\log\frac{K}{\delta}\alpha^2 - C_2\max_k |f^{(p)}(\boldsymbol{\mu}_k;\boldsymbol{\theta}(t)|),$
(E.15)

Whenever $c_{kj} \ge \sqrt{\frac{4}{5}}$,

$$\frac{d}{dt}c_{kj} \ge pc_{kj}^{p-1}\left(1 - \frac{1}{2^{p-2}}\right)(1 - c_{kj}) - C_1\log\frac{K}{\delta}\alpha^2 - C_2\max_k|f^{(p)}(\boldsymbol{\mu}_k;\boldsymbol{\theta}(t))|, \quad (E.16)$$

1744 Notice that by Lemma 8 and (E.12), we have

$$\max_{i} |f^{(p)}(\boldsymbol{x}_{i};\boldsymbol{\theta})| \leq \frac{p\tilde{\Delta}_{1}^{p-1}\tilde{\Delta}_{2}\alpha^{2}}{4}$$
(E.17)

These suffices to show that c_{kj} will reach $1 - \frac{C\alpha^2}{2}$ in less than \bar{t}_1 time.

For some choice of C and sufficiently small α , we have: Whenever, $\tilde{\Delta}_1 \leq c_{kj} \leq \sqrt{\frac{4}{5}}$, 1750 1751 $\frac{d}{dt}c_{kj} \ge pc_{kj}^{p-1}\tilde{\Delta}_2(1-c_{kj}) - C_1\log\frac{K}{\delta}\alpha^2 - C_2\max_{h}|f^{(p)}(\boldsymbol{\mu}_k;\boldsymbol{\theta}(t))|$ 1752 1753 $\geq p c_{kj}^{p-1} \tilde{\Delta}_2 \left(1 - \sqrt{\frac{4}{5}} \right) - C_1 \log \frac{K}{\delta} \alpha^2 - C_2 \max_k |f^{(p)}(\boldsymbol{\mu}_k; \boldsymbol{\theta}(t))|$ 1754 1755 1756 $\geq p c_{kj}^{p-1} \tilde{\Delta}_2 \left(1 - \sqrt{\frac{4}{5}} \right) - C_1 \log \frac{K}{\delta} \alpha^2 - C_2 \frac{p \tilde{\Delta}_1^{p-1} \tilde{\Delta}_2 \alpha^2}{4} \,.$ 1757 1758 1759 $\geq \frac{p}{2}c_{kj}^{p-1}\tilde{\Delta}_2\left(1-\sqrt{\frac{4}{5}}\right) \geq \frac{p}{20}c_{kj}^{p-1}\tilde{\Delta}_2\,,$ 1760 (E.18) 1761

where we uses the fact that $c_{kj} \ge \tilde{\Delta}_1$ in the last inequality. Whenever, $\sqrt{\frac{4}{5}} \le c_{kj} \le 1 - \frac{C\alpha^2}{2}$, $d \ge p-1\left(1, \frac{1}{2}\right)(1, \dots) = C_1 \frac{K}{2} + C_2 \frac{1}{2} + C_2 \frac{1}{2}$

$$\frac{d}{dt}c_{kj} \ge pc_{kj}^{p-1} \left(1 - \frac{1}{2^{p-2}}\right) (1 - c_{kj}) - C_1 \log \frac{K}{\delta} \alpha^2 - C_2 \max_k |f^{(p)}(\boldsymbol{\mu}_k; \boldsymbol{\theta}(t)|)|$$
$$\ge p(2^{p-1} - 2)(1 - c_{kj}) - C_1 \log \frac{K}{\delta} \alpha^2 - C_2 \max_k |f^{(p)}(\boldsymbol{\mu}_k; \boldsymbol{\theta}(t)|)|$$

1766 1767 1768

$$\geq p(2^{p-1}-2)(1-c_{kj}) - C_1 \log \frac{\pi}{\delta} \alpha^2 - C_2 \frac{p-1}{4} \frac{22\alpha}{4}$$
$$\geq \frac{p}{2}(2^{p-1}-2)(1-c_{kj}), \qquad (E.19)$$

1773 1774 where we uses the fact that $c_{kj} \leq 1 - C \log \frac{K}{\delta} \alpha^2$ in the last inequality. The right-hand sides of 1775 (E.18) and (E.19) is positive, which proves that c_{kj} is monotonically increasing before reaching 1776 $1 - C \log \frac{K}{\delta} \alpha^2$. Lastly,

1777

1778 1. by Lemma 9 and (E.18), it takes at most $\frac{20}{(p-2)p\tilde{\Delta}_2\tilde{\Delta}_1^{p-2}}$ time for c_{kj} to travel from $\tilde{\Delta}_1$ to $\sqrt{\frac{4}{5}}$; 1779 1780 2. by Lemma 10 and (E.10), it takes at most $\frac{2}{(p-2)p\tilde{\Delta}_2\tilde{\Delta}_1^{p-2}}$ large $\frac{1}{p}$ time for c_{kj} to travel from $\sqrt{\frac{4}{5}}$;

2. by Lemma 10 and (E.19), it takes at most
$$\frac{2}{p(2^{p-1}-2)}\log\frac{1}{C\log\frac{K}{\delta}\alpha^2}$$
 time for c_{kj} to travel from $\sqrt{\frac{4}{5}}$ to $1 - C\log\frac{K}{\delta}\alpha^2$.

Therefore, we have 1783

$$t_{1j}^{(k)} := \inf\{t : c_{kj}(t) \ge 1 - \frac{\alpha^2}{2}\} \le \frac{20}{(p-2)p\tilde{\Delta}_2\tilde{\Delta}_1^{p-2}} + \frac{2}{p(2^{p-1}-2)}\log\frac{1}{C\log\frac{K}{\delta}\alpha^2} = \bar{t}_1,$$
(E.20)

which contradicts our initial assumption that $c_{kj}(t) > \bar{t}_1$. Hence $t_1 \le \bar{t}_1$.

Maintaining $C \log \frac{K}{\delta} \alpha^2$ alignment until t_2 We have shown that at some $t_1 \leq \bar{t}_1$, all c_{kj} have grown above $1 - C \log \frac{K}{\delta} \alpha^2$. Now we show that any $c_{kj}(t)$ stays above $1 - C \log \frac{K}{\delta} \alpha^2$ between $[t_1, t_2]$. It suffices to show that for any $t \leq t_2$,

$$\left. \frac{d}{dt} c_{kj} \right|_{c_{kj} = 1 - C \log \frac{K}{\delta} \alpha^2} \ge 0.$$
(E.21)

1795 Indeed, the inequality (E.16) is still valid before t_2 , i.e.

$$\frac{d}{dt}c_{kj} \ge pc_{kj}^{p-1}\left(1 - \frac{1}{2^{p-2}}\right)(1 - c_{kj}) - C_1\log\frac{K}{\delta}\alpha^2 - C_2\max_k |f^{(p)}(\boldsymbol{\mu}_k;\boldsymbol{\theta}(t)|)|$$
$$\ge p(2^{p-1} - 2)(1 - c_{kj}) - C_1\log\frac{K}{\delta}\alpha^2 - C_2\frac{p\tilde{\Delta}_1^{p-1}\tilde{\Delta}_2\alpha^2}{4}.$$

1801 Therefore, for some choice of C and sufficiently small α ,

$$\frac{d}{dt}c_{kj}\Big|_{c_{kj}=1-C\log\frac{K}{\delta}\alpha^2} \ge p(2^{p-1}-2)C\log\frac{K}{\delta}\alpha^2 - C_1\log\frac{K}{\delta}\alpha^2 - C_2\frac{p\tilde{\Delta}_1^{p-1}\tilde{\Delta}_2\alpha^2}{4} \ge 0.$$
(E.22)

Hence

$$\min_{k} \min_{j \in \mathcal{N}_{k}} c_{kj}(t) \ge 1 - C \log \frac{K}{\delta} \alpha^{2}, \forall t \in [t_{1}, t_{2}].$$
(E.23)

Г		1
L		L 1
5	-	

1836
1837
1838FPRELU CONVERGES TO OPTIMAL ℓ_2 -ROBUST CLASSIFIER, PART FOUR:
CONVERGENCE PHASE1838CONVERGENCE PHASE

1839 1840 F.1 AXUILIARY LEMMAS

¹⁸⁴¹ We need the following lemmas (proofs provided in Appendix G):

Lemma 11. Let p > 2. Condition on good event \mathcal{E}_{good} . Suppose the following is true at some point on the GF trajectory:

1845 *I.* $c_{kj}(t) \ge 1 - 2C_a \log \frac{K}{\delta} \alpha^2, \forall k, j \in \mathcal{N}_k;$

1846 1847 2. $\sum_{j \in \mathcal{N}_k} \|\boldsymbol{w}_j\|^2 \le 1 + C_w \log \frac{K}{\delta} \alpha^2, \ \forall k;$

1848 1849 3. $\sum_{j \in \mathcal{N}_c} \| \boldsymbol{w}_j \|^2 = \tilde{o}(\alpha^2).$

1850 Then the following holds for every $1 \le k \le K$, $i \in \mathcal{I}_k$, 1851

$$f^{(p)}(\boldsymbol{x}_i; \boldsymbol{\theta}) \leq \sum_{j \in \mathcal{N}_k} \|\boldsymbol{w}_j\|^2 \left(1 + 2^{p+2}C\sqrt{\log\frac{K^2N}{\delta}}\alpha^2 \right) + 2KC\alpha^p;$$

$$f^{(p)}(\boldsymbol{x}_i; \boldsymbol{\theta}) \geq \sum_{j \in \mathcal{N}_k} \|\boldsymbol{w}_j\|^2 \left(1 - 4pC\sqrt{\log\frac{K^2N}{\delta}}\alpha^2 \right) - 2KC\alpha^p.$$

1855 1856 1857

1867 1868

1870 1871

1880

1852 1853 1854

Lemma 12. Let p > 2. Condition on good event \mathcal{E}_{good} . Suppose the following is true at some point on the GF trajectory:

1860 I. $c_{kj}(t) \ge 1 - 2C_a \log \frac{K}{\delta} \alpha^2, \ \forall k, j \in \mathcal{N}_k;$ 1861

1862 2.
$$\sum_{i \in N_{L}} \| \boldsymbol{w}_{i} \|^{2} \leq 1 + C_{w} \log \frac{K}{\delta} \alpha^{2}, \forall k;$$

1863 1864 3. $\sum_{j \in \mathcal{N}_c} \| \boldsymbol{w}_j \|^2 = \tilde{o}(\alpha^2).$

1865 1866 Furthermore, suppose additionally that for some $k, j \in \mathcal{N}_k$:

$$1 - 2C_a \log \frac{K}{\delta} \alpha^2 \le c_{kj}(t) \le 1 - C_a \log \frac{K}{\delta} \alpha^2;$$

1869 Then the following holds for the same k, j,

$$\frac{d}{dt}c_{kj} \ge -CK\log\frac{K^2N}{\delta}\alpha^{\min\{p,4\}}$$

Lemma 13. Let p > 2. Condition on good event \mathcal{E}_{good} . Suppose the following is true at some point on the GF trajectory :

$$\begin{array}{ll} & \text{1875} & I. \ c_{kj}(t) \geq 1 - 2C_a \log \frac{K}{\delta} \alpha^2, \ k, j \in \mathcal{N}_k; \\ & \text{1876} & \\ & \text{1877} & 2. \ \sum_{j \in \mathcal{N}_k} \| \boldsymbol{w}_j \|^2 \leq 1 + C_w \log \frac{K}{\delta} \alpha^2, \ \forall k; \\ & \text{1878} & \\ & \text{1879} & 3. \ \sum_{j \in \mathcal{N}_c} \| \boldsymbol{w}_j \|^2 = \tilde{o}(\alpha^2). \end{array}$$

Then the following holds for every $1 \le k \le K$,

$$\frac{d}{dt}\left(\sum_{j\in\mathcal{N}_k}\|\boldsymbol{w}_j\|^2\right) \leq 2\left(1-\sum_{j\in\mathcal{N}_k}\|\boldsymbol{w}_j\|^2+C\log\frac{K}{\delta}\alpha^2\right)\left(\sum_{j\in\mathcal{N}_k}\|\boldsymbol{w}_j\|^2\right)\,,$$

1885 and

1886
1887
1888
1889
1889
1889
1889
1890

$$\frac{d}{dt}\left(\sum_{j\in\mathcal{N}_k}\|\boldsymbol{w}_j\|^2\right) \ge 2\left(1-\sum_{j\in\mathcal{N}_k}\|\boldsymbol{w}_j\|^2-C\log\frac{K}{\delta}\alpha^2\right)\left(\sum_{j\in\mathcal{N}_k}\|\boldsymbol{w}_j\|^2\right),$$
1886
1887
1888
1889
1890
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900
1900

where C is some universal constant such that $C < C_w$.

Lemma 14. Consider the same assumptions as in Proposition 2. Given the t_1 in Proposition 2, the following holds $\forall 1 \le k \le K$:

1895

1897 1898

1902

1905 1906 1907

1908 1909

1910 1911 1912

$$\sum_{j \in \mathcal{N}_k} \|\boldsymbol{w}_j(t_1)\|^2 \ge \exp\left(-\frac{2p^{p+2}K}{p(p-2)\tilde{\Delta}_2\tilde{\Delta}_1^{p-2}}\right) W_{\min}^2 \epsilon^2.$$
(F.1)

Lemma 15. Given some $0 < \Delta < \frac{1}{4}$, if for some z(t), the following holds

$$\frac{d}{dt}z \ge (1 - z - \Delta)z, \ z(0) = z_0, \ z(T) = z_1,$$
(F.2)

for some $0 < z_0 \leq \frac{1}{4}$, and $z_0 \leq z_1 < 1 - \Delta$. Then the travel time T for z(t) to go from z_0 to z_1 satisfies:

$$T \le 2\left(\log\frac{1}{1-z_1-\Delta} + \log\frac{1}{z_0}\right). \tag{F.3}$$

Lemma 16. Condition on good event \mathcal{E}_{good} , we have

$$\sum_{j \in \mathcal{N}_c} \|\boldsymbol{w}_j(t)\|^2 = \tilde{o}(\alpha^2), \ \forall t \le T^* \,.$$
(F.4)

Lemma 17. If the neurons $\{w_j\}_{j=1}^h$ satisfies the following for some $0 \le \delta \le 1$ and $\nu, \zeta > 0$:

• $\max_k \max_{j \in \mathcal{N}_k} c_{kj}(t) \ge 1 - \delta;$

•
$$\left|1-\sum_{j\in\mathcal{N}_k}\|\boldsymbol{w}_j\|^2\right|\leq\nu,$$

1913 • $\sum_{j\in\mathcal{N}^c}\|w_j\|^2\leq \zeta$,

1915 then $\sup_{\boldsymbol{x}\in\mathbb{S}^{D-1}} \left| f^{(p)}(\boldsymbol{x};\boldsymbol{\theta}) - F^{(p)}(\boldsymbol{x}) \right| \le K(1+\nu)(2^p-1)2\delta + K\nu + \zeta$

1917 F.2 Proof of Theorem 3

1919 **Theorem 3** (Restated). Let p > 2. Given $0 \le \delta \le 1$ and a sufficiently small α_0^2 , consider data dimension $D \ge \tilde{\Omega}(\alpha_0^{-2})$ and per-cluster sample size $\tilde{\Omega}(\alpha_0^{-2}) \le N \le \tilde{o}(\exp(\alpha_0^{-2}))$. With probability 1920 1921 at least $1 - \delta$, the GF dynamics with a balanced dataset $\hat{D} = \{x_i, y_i\}_{i=1}^{KN}$ sampled with intra-cluster 1922 variance $\alpha^2 \leq \alpha_0^2$, starting from some ϵ -small and balanced (Assumption 1) initialization $\theta(0)$ that 1923 satisfies Assumption 2 with a non-degeneracy gap $\Delta = \Theta(1)$ and has a sufficiently small initialization 1924 scale $\epsilon = \tilde{\Theta}(\alpha_0^{8K})$, leads to a solution $\theta(t), t \geq 0$ such that: for some $t^* = \tilde{O}(\log \frac{1}{\alpha_0})$ and 1925 $T^* = \tilde{\Theta}\left(\log \frac{1}{\alpha_0}\right) + \tilde{\Omega}\left(\frac{1}{\alpha_0^{\min\{p-2,2\}}}\right) \text{ with } [t^*, T^*] \neq \emptyset, \text{ we have } \mathcal{L}(\boldsymbol{\theta}(t)) = \tilde{\mathcal{O}}(\alpha_0^4), \forall t \in [t^*, T^*]$ 1926 1927 and $\sup_{t \in [t^*, T^*]} \sup_{\boldsymbol{x} \in \mathbb{S}^{D-1}} \left| f^{(p)}(\boldsymbol{x}; \boldsymbol{\theta}(t)) - F^{(p)}(\boldsymbol{x}) \right| \leq \tilde{\mathcal{O}}\left(\alpha_0^2\right) \,.$ 1928 (F.5) 1929

1931 *Proof.* We have shown in Proposition 2, and Lemma 14 that:

1933
1. Any
$$c_{kj}(t)$$
 staying above $1 - C_a \log \frac{K}{\delta} \alpha^2$ during $[t_1, t_2]$;
1935
2. $\sum_{j \in \mathcal{N}_k} \|\boldsymbol{w}_j(t_1)\|^2 \ge \exp\left(-\frac{2p^{p+2}K}{p(p-2)\tilde{\Delta}_2\tilde{\Delta}_1^{p-2}}\right) W_{\min}^2 \epsilon^2$, for every $1 \le 1$

We define

$$t^{*} = \underbrace{t_{1}}_{\mathcal{O}(1)} + 2\left(\log\frac{1}{(C - C_{w})\log\frac{K}{\delta}\alpha^{2}} + \frac{2p^{p+2}K}{p(p-2)\tilde{\Delta}_{2}\tilde{\Delta}_{1}^{p-2}} + \log\frac{1}{W_{\min}^{2}\epsilon^{2}}\right)$$
(F.6)

 $k \leq K$.

1941 1942 1943

1937

1938 1939

$$T^* = \underbrace{t_2}_{\Theta(\log \frac{1}{\epsilon})} + \frac{C}{\log \frac{K^2 N}{\delta}} \alpha^{\max\{2-p,-2\}}$$
(F.7)

1944 Since $\epsilon = \Theta(\alpha^{8K})$. For sufficiently small α , we have $\mathcal{O}(\log \frac{1}{\alpha}) = t^* \leq T^* = \Theta(\alpha^{\min\{2-p,-2\}})$. 1945 Our goal is to show that 1946 1947 1. Before T^* , one must have $\max_k \max_{j \in \mathcal{N}_k} c_{kj}(t) \ge 1 - 2C_a \log \frac{K}{\delta} \alpha^2$ and $\sum_{j \in \mathcal{N}_k} \|\boldsymbol{w}_j(t)\|^2 \le 1 - 2C_a \log \frac{K}{\delta} \alpha^2$ 1948 $1 + C_w \log \frac{K}{\delta} \alpha^2;$ 1949 2. Before T^* , for all k, whenever $\sum_{i \in \mathcal{N}_k} \| \boldsymbol{w}_i(t) \|^2$ reaches $1 - C_w \log \frac{K}{\delta} \alpha^2$, it can not drop below 1950 1951 $1 - C_w \log \frac{K}{\delta} \alpha^2;$ 1952 3. After t^* , for all k, one must have $\sum_{i \in \mathcal{N}_k} \|\boldsymbol{w}_i\|^2 \ge 1 - 2C_w \log \frac{K}{\delta} \alpha^2$. 1953 1954 We also have $\sum_{j \in \mathcal{N}_c} \|\boldsymbol{w}_j\|^2 = \tilde{o}(\alpha^2)$, then applying Lemma 17 gives the desired result. The 1955 statement that $\mathcal{L}(t) = \tilde{O}(\alpha^4)$ is due to the fact that $|y_i - f^{(p)}(\boldsymbol{x}_i; \boldsymbol{\theta}(t))| = \tilde{O}(\alpha^2)$ during $[t^*, T^*]$. 1957 **First claim**: The two inequalities hold before t_2 , thus it suffices to study 1958 1959 $\tau_3 := \inf \left\{ t \ge t_2 : \max_{k} \max_{j \in \mathcal{N}_k} c_{kj}(t) \le 1 - 2C \log \frac{K}{\delta} \alpha^2 \right\} ,$ 1960 1961 $\tau_4 := \inf \left\{ t \ge t_2 : \sum_{i \in \mathcal{N}_i} \|\boldsymbol{w}_j\|^2 \ge 1 + C \log \frac{K}{\delta} \alpha^2 \right\} ,$ 1963 1964 and show that $\min\{\tau_3, \tau_4\} \ge T^*$. We proof it by contradiction, suppose $\min\{\tau_3, \tau_4\} \le T^*$, then it 1965 must be either $\tau_3 = \min\{\tau_3, \tau_4\} \le T^*$ or $\tau_4 = \min\{\tau_3, \tau_4\} \le T^*$. 1966 Consider the first case that $\tau_3 = \min\{\tau_3, \tau_4\} \leq T^*$, then there exists some k and $j \in \mathcal{N}_k$ and some 1967 $\tau_{3^{-}} \geq t_2$ such that 1968 $1 - 2C\log\frac{K}{s}\alpha^2 \le c_{kj}(t) \le 1 - C\log\frac{K}{s}\alpha^2, \forall t \in [\tau_{3^-}, \tau_3],$ 1969 (F.8) 1970 $c_{kj}(\tau_{3^{-}}) = 1 - C \log \frac{K}{\delta} \alpha^2, \ c_{kj}(\tau_3) = 1 - 2C \log \frac{K}{\delta} \alpha^2$ 1971 (F.9) 1972 since $c_{kj}(t)$ is continuous and has to travel from $1 - C \log \frac{K}{\delta} \alpha^2$ to $1 - 2C \log \frac{K}{\delta} \alpha^2$. By Lemma 6, 1973 we have 1974 $\frac{d}{dt}c_{kj} \ge -CK\log\frac{K^2N}{\delta}\alpha^{\min\{p,4\}}, \forall t \in [\tau_{3^-}, \tau_3].$ 1975 1976 Then by the fundamental theorem of calculus, we have 1977 1978 $-C\log\frac{K}{\delta}\alpha^{2} = c_{kj}(\tau_{3}) - c_{kj}(\tau_{3^{-}}) = \int_{\tau_{a^{-}}}^{\tau_{3}} \frac{d}{dt}c_{kj} \ge \int_{\tau_{a^{-}}}^{\tau_{3}} -CK\log\frac{K^{2}N}{\delta}\alpha^{\min\{p,4\}}$ 1979 $= -(\tau_3 - \tau_{3^-})CK\log \frac{K^2N}{\delta} \alpha^{\min\{p,4\}}$ 1981 1982 (F.10) 1983 Therefore, for some constant C > 0, 1984 $(\tau_3 - \tau_{3^-}) \geq \frac{C}{\log \frac{K^2 N}{s}} \alpha^{\max\{2-p,-2\}} \Rightarrow (\tau_3 - t_2) \geq \frac{C}{\log \frac{K^2 N}{\delta}} \alpha^{\max\{2-p,-2\}} \,.$ 1985 (F.11) 1986

 $[t_2, \tau_3]$ has length at least $\frac{C}{\log \frac{K^2 N}{\delta}} \alpha^{\max\{2-p, -2\}}$ thus is an interval that contains $[t_2, T^*]$. Contradict-1987 1988 ing our assumption that $\tau_3 \stackrel{\circ}{\leq} \stackrel{\circ}{T^*}$. The case one is thus eliminated. 1989

Consider the second case that $\tau_4 = \min\{\tau_3, \tau_4\} \le T^*$, then by the continuity of $||w_j||$, we know that 1990 there exists some k such that $\sum_{j \in \mathcal{N}_k} \| \boldsymbol{w}_j(\tau_4) \|^2 = 1 + C_w \log \frac{K}{\delta} \alpha^2$. However, by Lemma 11, we 1992 have, at τ_4 , 1 \

1994
1995
1996
1997
$$\frac{d}{dt}\left(\sum_{j\in\mathcal{N}_k}\|\boldsymbol{w}_j\|^2\right) \le 2\left(1 - \sum_{\substack{j\in\mathcal{N}_k\\ =1+C_w\log\frac{K}{\delta}\alpha^2}}\|\boldsymbol{w}_j\|^2 + C\log\frac{K}{\delta}\alpha^2\right)\left(\sum_{j\in\mathcal{N}_k}\|\boldsymbol{w}_j\|^2\right)$$

1998
1999
2000
2001
$$= 2\left((C-C_w)\log\frac{K}{\delta}\alpha^2\right)\left(\sum_{j\in\mathcal{N}_k}\|\boldsymbol{w}_j\|^2\right) < 0,$$

which indicates that $\sum_{j \in \mathcal{N}_k} \| \boldsymbol{w}_j \|^2$ can not surpass $1 + C_w \log \frac{K}{\delta} \alpha^2$ after τ_4 , violating the definition of τ_4 , leading to a contradiction. Therefore the second case is eliminated as well. We must have $\min\{\tau_3, \tau_4\} \ge T^*$. The first claim is proved.

Second claim By Lemma 13 (it applies to any $t \leq T^*$ given the proof in our first step), we have

$$\frac{d}{dt} \left(\sum_{j \in \mathcal{N}_k} \|\boldsymbol{w}_j\|^2 \right) \bigg|_{\sum_{j \in \mathcal{N}_k} \|\boldsymbol{w}_j\|^2 = 1 - C_w \log \frac{K}{\delta} \alpha^2} \ge 2 \left(C_w \log \frac{K}{\delta} \alpha^2 - C \log \frac{K}{\delta} \alpha^2 \right) \left(\sum_{j \in \mathcal{N}_k} \|\boldsymbol{w}_j\|^2 \right) + 2 O .$$

Therefore, whenever $\sum_{j \in \mathcal{N}_k} \| \boldsymbol{w}_j(t) \|^2$ reaches $1 - C_w \log \frac{K}{\delta} \alpha^2$, it can not drop below $1 - C_w \log \frac{K}{\delta} \alpha^2$ $C_w \log \frac{K}{\delta} \alpha^2$. The second claim is proved.

Third claim Lastly, we just need an upper bound on the travel time for $\sum_{j \in N_k} \| \boldsymbol{w}_j(t) \|^2$ to go from $\sum_{j \in \mathcal{N}_k} \| \boldsymbol{w}_j(t_1) \|^2$ to $1 - C_w \log \frac{K}{\delta} \alpha^2$, for which we simply combine Lemma 13, 14, and 15 to see the travel time is upper bounded by

$$2\left(\log\frac{1}{(C-C_w)\log\frac{K}{\delta}\alpha^2} + \frac{2p^{p+2}K}{p(p-2)\tilde{\Delta}_2\tilde{\Delta}_1^{p-2}} + \log\frac{1}{W_{\min}^2\epsilon^2}\right).$$
 (F.12)

Thus
$$\sum_{j \in \mathcal{N}_k} \| \boldsymbol{w}_j(t) \|^2$$
 must reach $1 - C_w \log \frac{K}{\delta} \alpha^2$ by t^* .

$\begin{array}{ccc} & & & \\ \textbf{C} & & \\ \textbf{C}$

Lemma 5 (Restated). *Given an initialization shape that satisfies Assumption 2 with non-degeneracy* gap $\Delta > 0$, then for $j \in \mathcal{N}_k$, we have

$$c_{kj}(0) = \cos(\boldsymbol{\mu}_k, \boldsymbol{w}_j(0)) \ge \sqrt{\frac{1}{2} \left(\frac{1}{(1-\Delta)^2} - 1\right)} := \tilde{\Delta}_1, \tag{G.1}$$

$$\frac{c_{lj}^{p-2}(0)}{c_{kj}^{p-2}(0)} \le (1 - \sqrt{2\Delta})^{p-2} := 1 - \tilde{\Delta}_2, \forall l \neq k \text{ with } y_l = y_k \text{ and } c_{lj}(0) > 0$$
(G.2)

Proof. We prove both inequalities by contradiction.

First inequality Suppose $0 < c_{kj}(0) = \cos(\boldsymbol{\mu}_k, \boldsymbol{w}_j(0)) = \cos(\boldsymbol{\mu}_k, \boldsymbol{w}_{j0}) < \tilde{\Delta}_1$, then consider $\tilde{\boldsymbol{w}}_{j0} = \frac{\boldsymbol{w}_{j0}}{\|\boldsymbol{w}_{j0}\|}$, and

$$\tilde{\boldsymbol{w}} = \tilde{\boldsymbol{w}}_{j0} - \frac{c_{kj}(0)}{1 - c_{kj}(0)} (\boldsymbol{\mu}_k - \tilde{\boldsymbol{w}}_{j0}).$$
(G.3)

Notice that here $c_{kj}(0) = \cos(\mu_k, w_{j0}) = \langle \mu_k, \tilde{w}_{j0} \rangle$. It is easy to verify that $\langle \mu_l, \tilde{w} \rangle = 0, \forall 1 \leq l \leq K$, thus $\tilde{w} \in \partial(\bigcup_{k \in \mathcal{K}} \mathcal{R}_k)$, and

$$d\left(\boldsymbol{w}_{j0},\partial\left(\bigcup_{k\in\mathcal{K}}\mathcal{R}_{k}\right)\right)=1-\sup_{\boldsymbol{w}\in\partial\left(\bigcup_{k\in\mathcal{K}}\mathcal{R}_{k}\right)}\cos\left(\tilde{\boldsymbol{w}}_{j0},\boldsymbol{w}\right)\leq1-\cos\left(\tilde{\boldsymbol{w}}_{j0},\tilde{\boldsymbol{w}}\right),\qquad(G.4)$$

Since one can compute

$$\cos(\tilde{\boldsymbol{w}}_{j0}, \tilde{\boldsymbol{w}}) = \frac{\langle \tilde{\boldsymbol{w}}_{j0}, \tilde{\boldsymbol{w}} \rangle}{\|\tilde{\boldsymbol{w}}_{j0}\| \|\tilde{\boldsymbol{w}}\|} = \frac{1 + c_{kj}(0)(1 - c_{kj}(0))}{\sqrt{1 + 2c_{kj}^2(0)}} \ge \frac{1}{\sqrt{1 + 2c_{kj}^2(0)}} > 1 - \Delta, \quad (G.5)$$

where the last inequality is due to our assumption that $c_{kj}(0) < \tilde{\Delta}_1$. Combining (G.4)(G.5), we have

$$d\left(\boldsymbol{w}_{j0}, \partial\left(\bigcup_{k\in\mathcal{K}}\mathcal{R}_{k}\right)\right) < \Delta, \qquad (G.6)$$

which contradicts our assumption that the non-degeneracy gap is at least Δ .

2088 Second inequality Suppose there exists an $l \neq k$ such that $y_l = y_k$ and $\frac{c_{l_j}^{p-2}(0)}{c_{k_j}^{p-2}(0)} > (1 - \sqrt{2\Delta})^{p-2}$ 2089 and $c_{l_j}(0) > 0$, we pick the *l* that has the largest $c_{l_j}(0)$, then consider $\tilde{w}_{j0} = \frac{w_{j0}}{\|w_{j0}\|}$, and

$$\tilde{w} = \tilde{w}_{j0} - \frac{c_{kj}(0) - c_{lj}(0)}{2} (\mu_k - \mu_l).$$
(G.7)

It can be verified that $\|\tilde{w}\| = 1$, $\cos(\mu_k, \tilde{w}) = \cos(\mu_l, \tilde{w}) = \frac{c_{kj}(0) + c_{lj}(0)}{2}$, and $\cos(\mu_m, \tilde{w}) = \frac{\cos(\mu_m, \tilde{w})}{2} = \frac{\cos(\mu_m, \tilde{w}_{j0})}{2} \le \cos(\mu_l, \tilde{w}), \forall m \neq k \text{ or } l$. All of the above together implies $\tilde{w} \in (\partial \mathcal{R}_k) \cap (\partial \mathcal{R}_l) \subset \partial(\bigcup_{k \in \mathcal{K}} \mathcal{R}_k)$, and

$$d\left(\boldsymbol{w}_{j0},\partial\big(\bigcup_{k\in\mathcal{K}}\mathcal{R}_k\big)\right) = 1 - \sup_{\boldsymbol{w}\in\partial\big(\bigcup_{k\in\mathcal{K}}\mathcal{R}_k\big)}\cos\left(\tilde{\boldsymbol{w}}_{j0},\boldsymbol{w}\right) \le 1 - \cos(\tilde{\boldsymbol{w}}_{j0},\tilde{\boldsymbol{w}}), \quad (G.8)$$

2100 One can compute

$$\cos(\tilde{\boldsymbol{w}}_{j0}, \tilde{\boldsymbol{w}}) = \frac{\langle \tilde{\boldsymbol{w}}_{j0}, \tilde{\boldsymbol{w}} \rangle}{\|\tilde{\boldsymbol{w}}_{i0}\| \|\tilde{\boldsymbol{w}}\|} = 1 - \frac{(c_{kj}(0) - c_{lj}(0))^2}{2}$$

2104
2105
$$\ge 1 - \frac{\left(1 - \frac{c_{lj}(0)}{c_{kj}(0)}\right)}{2}$$

$$\begin{array}{l} 2106\\ 2107\\ 2108\\ 2109\\ 2110\\ 2111\\ 2112 \end{array} \\ \geq 1 - \frac{\left(1 - \left(\frac{c_{l_j}^{p-2}(0)}{c_{k_j}^{p-2}(0)}\right)^{\frac{1}{p-2}}\right)^2}{2}\\ \geq 1 - \frac{\left(1 - (1 - \tilde{\Delta}_2)^{\frac{1}{p-2}}\right)^2}{2} = 1 - \Delta, \end{array}$$
(G.9)

where the last inequality is due to our assumption that $\frac{c_{lj}^{p-2}(0)}{c_{kj}^{p-2}(0)} > (1 - \sqrt{2\Delta})^{p-2}$. Combining (G.4)(G.5), we have

$$d\left(\boldsymbol{w}_{j0},\partial\left(\bigcup_{k\in\mathcal{K}}\mathcal{R}_{k}\right)\right)<\Delta,$$
(G.10)
that the non-degeneracy gap is at least Δ .

which contradicts our assumption that the non-degeneracy gap is at least Δ .

Lemma 6 (Restated). Let p > 2. Condition on good event \mathcal{E}_{good} . Given some $1 \le k \le K$ and some $j \in \mathcal{N}_k$ and suppose the following is true at some point on the GF trajectory:

$$\begin{array}{ll} \mathbf{2123} \\ \mathbf{2124} \end{array} \quad l. \ c_{kj} \geq \tilde{\Delta}_{1j} \end{array}$$

2125 2.
$$\frac{|c_{lj}|}{c_{kj}} \le (1 - \sqrt{2\Delta}), \forall l \neq k.$$

2126

Then the following holds:

$$\frac{d}{dt}c_{kj} \ge pc_{kj}^{p-1}\tilde{\Delta}_2(1-c_{kj}) - C_1\log\frac{K}{\delta}\alpha^2 - C_2\max_k |f^{(p)}(\boldsymbol{\mu}_k;\boldsymbol{\theta}(t))|,$$

for some universal constant C_1, C_2 that depends on p. If one further assume $c_{kj} \ge \sqrt{\frac{4}{5}}$, then the lower bound can be improved as

$$\frac{d}{dt}c_{kj} \ge pc_{kj}^{p-1}\left(1 - \frac{1}{2^{p-2}}\right)\left(1 - c_{kj}\right) - C_1\log\frac{K}{\delta}\alpha^2 - C_2\max_i |f^{(p)}(\boldsymbol{x}_i;\boldsymbol{\theta})|,$$

Proof. When $1 \le k \le K_1, j \in \mathcal{N}_k$ implies that $j \in \mathcal{N}_+$ thus $sign(v_j) = 1$. We shall primarily focus on this case as the proof is nearly identical for $K_1 + 1 \le k \le K$.

$$\frac{d}{dt}c_{kj} = -\frac{1}{N}\sum_{i:\langle \boldsymbol{x}_{i}, \boldsymbol{w}_{j} \rangle > 0} \nabla_{\hat{y}}\ell_{i} p\left(\left\langle \boldsymbol{x}_{i}, \frac{\boldsymbol{w}_{j}}{\|\boldsymbol{w}_{j}\|}\right\rangle\right)^{p-1}\left(\left\langle \boldsymbol{\mu}_{k}, \boldsymbol{x}_{i}\right\rangle - \left\langle \boldsymbol{x}_{i}, \frac{\boldsymbol{w}_{j}}{\|\boldsymbol{w}_{j}\|}\right\rangle c_{kj}\right) \\
= \frac{1}{N}\sum_{i:\langle \boldsymbol{x}_{i}, \boldsymbol{w}_{j} \rangle > 0} (y_{i} - f^{(p)}(\boldsymbol{x}_{i}; \boldsymbol{\theta})) p\left(\left\langle \boldsymbol{x}_{i}, \frac{\boldsymbol{w}_{j}}{\|\boldsymbol{w}_{j}\|}\right\rangle\right)^{p-1}\left(\left\langle \boldsymbol{\mu}_{k}, \boldsymbol{x}_{i}\right\rangle - \left\langle \boldsymbol{x}_{i}, \frac{\boldsymbol{w}_{j}}{\|\boldsymbol{w}_{j}\|}\right\rangle c_{kj}\right) \\
= \frac{1}{N}\sum_{i:\langle \boldsymbol{x}_{i}, \boldsymbol{w}_{j} \rangle > 0} y_{i} p\left(\left\langle \boldsymbol{x}_{i}, \frac{\boldsymbol{w}_{j}}{\|\boldsymbol{w}_{j}\|}\right\rangle\right)^{p-1}\left(\left\langle \boldsymbol{\mu}_{k}, \boldsymbol{x}_{i}\right\rangle - \left\langle \boldsymbol{x}_{i}, \frac{\boldsymbol{w}_{j}}{\|\boldsymbol{w}_{j}\|}\right\rangle c_{kj}\right) \\
= \frac{1}{N}\sum_{i:\langle \boldsymbol{x}_{i}, \boldsymbol{w}_{j} \rangle > 0} f^{(p)}(\boldsymbol{x}_{i}; \boldsymbol{\theta}) p\left(\left\langle \boldsymbol{x}_{i}, \frac{\boldsymbol{w}_{j}}{\|\boldsymbol{w}_{j}\|}\right\rangle\right)^{p-1}\left(\left\langle \boldsymbol{\mu}_{k}, \boldsymbol{x}_{i}\right\rangle - \left\langle \boldsymbol{x}_{i}, \frac{\boldsymbol{w}_{j}}{\|\boldsymbol{w}_{j}\|}\right\rangle c_{kj}\right) \\
= \frac{1}{N}\sum_{i:\langle \boldsymbol{x}_{i}, \boldsymbol{w}_{j} \rangle > 0} f^{(p)}(\boldsymbol{x}_{i}; \boldsymbol{\theta}) p\left(\left\langle \boldsymbol{x}_{i}, \frac{\boldsymbol{w}_{j}}{\|\boldsymbol{w}_{j}\|}\right\rangle\right)^{p-1}\left(\left\langle \boldsymbol{\mu}_{k}, \boldsymbol{x}_{i}\right\rangle - \left\langle \boldsymbol{x}_{i}, \frac{\boldsymbol{w}_{j}}{\|\boldsymbol{w}_{j}\|}\right\rangle c_{kj}\right) \\
= \frac{1}{N}\sum_{i:\langle \boldsymbol{x}_{i}, \boldsymbol{w}_{j} \rangle > 0} y_{i} p\left(\left\langle \boldsymbol{x}_{i}, \frac{\boldsymbol{w}_{j}}{\|\boldsymbol{w}_{j}\|}\right\rangle\right)^{p-1}\left(\left\langle \boldsymbol{\mu}_{k}, \boldsymbol{x}_{i}\right\rangle - \left\langle \boldsymbol{x}_{i}, \frac{\boldsymbol{w}_{j}}{\|\boldsymbol{w}_{j}\|}\right\rangle c_{kj}\right) \\
= \frac{1}{N}\sum_{i:\langle \boldsymbol{x}_{i}, \boldsymbol{w}_{j} \rangle > 0} y_{i} p\left(\left\langle \boldsymbol{x}_{i}, \frac{\boldsymbol{w}_{j}}{\|\boldsymbol{w}_{j}\|}\right\rangle\right)^{p-1}\left(\left\langle \boldsymbol{\mu}_{k}, \boldsymbol{x}_{i}\right\rangle - \left\langle \boldsymbol{x}_{i}, \frac{\boldsymbol{w}_{j}}{\|\boldsymbol{w}_{j}\|}\right\rangle c_{kj}\right) \\ = \frac{1}{N}\sum_{i\in\mathcal{I}_{k}:\langle \boldsymbol{x}_{i}, \boldsymbol{w}_{j} \rangle > 0} y_{i} p\left(\left\langle \boldsymbol{x}_{i}, \frac{\boldsymbol{w}_{j}}{\|\boldsymbol{w}_{j}\|}\right\rangle\right)^{p-1}\left(\left\langle \boldsymbol{\mu}_{k}, \boldsymbol{x}_{i}\right\rangle - \left\langle \boldsymbol{x}_{i}, \frac{\boldsymbol{w}_{j}}{\|\boldsymbol{w}_{j}\|}\right\rangle c_{kj}\right) \right)$$

(a)

$$-\underbrace{\frac{1}{N}\sum_{l\neq k}\sum_{i\in\mathcal{I}_{l}:\langle\boldsymbol{x}_{i},\boldsymbol{w}_{j}\rangle>0}y_{i}p\left(\left\langle\boldsymbol{x}_{i},\frac{\boldsymbol{w}_{j}}{\|\boldsymbol{w}_{j}\|}\right\rangle\right)^{p-1}\left(\left\langle\boldsymbol{\mu}_{k},\boldsymbol{x}_{i}\right\rangle-\left\langle\boldsymbol{x}_{i},\frac{\boldsymbol{w}_{j}}{\|\boldsymbol{w}_{j}\|}\right\rangle c_{kj}\right)}_{(b)}\right)+\Gamma_{1}$$
(G.11)

We handle these two terms differently:

$$(a) = \frac{1}{N} \sum_{i \in \mathcal{I}_k: \langle \boldsymbol{x}_i, \boldsymbol{w}_j \rangle > 0} y_i \, p\left(\left\langle \boldsymbol{x}_i, \frac{\boldsymbol{w}_j}{\|\boldsymbol{w}_j\|} \right\rangle\right)^{p-1} \left(\left\langle \boldsymbol{\mu}_k, \boldsymbol{x}_i \right\rangle - \left\langle \boldsymbol{x}_i, \frac{\boldsymbol{w}_j}{\|\boldsymbol{w}_j\|} \right\rangle c_{kj}\right)$$
$$= \frac{1}{N} \sum_{i \in \mathcal{I}_k} p\left(\left\langle \boldsymbol{\mu}_k + \boldsymbol{\varepsilon}_i, \frac{\boldsymbol{w}_j}{\|\boldsymbol{w}_j\|} \right\rangle\right)^{p-1} \left(\left\langle \boldsymbol{\mu}_k, \boldsymbol{\mu}_k + \boldsymbol{\varepsilon}_i \right\rangle - \left\langle \boldsymbol{\mu}_k + \boldsymbol{\varepsilon}_i, \frac{\boldsymbol{w}_j}{\|\boldsymbol{w}_j\|} \right\rangle c_{kj}\right)$$
$$= \frac{1}{N} \sum_{i \in \mathcal{I}_k} p\left(\left\langle \boldsymbol{\mu}_k + \boldsymbol{\varepsilon}_i, \frac{\boldsymbol{w}_j}{\|\boldsymbol{w}_j\|} \right\rangle\right)^{p-1} \left(1 - c_{kj}^2 + \left\langle \boldsymbol{\mu}_k, \boldsymbol{\varepsilon}_i \right\rangle - \left\langle \boldsymbol{\varepsilon}_i, \frac{\boldsymbol{w}_j}{\|\boldsymbol{w}_j\|} \right\rangle c_{kj}\right)$$
$$= \frac{1}{N} \sum_{i \in \mathcal{I}_k} p\left(\left\langle \boldsymbol{\mu}_k + \boldsymbol{\varepsilon}_i, \frac{\boldsymbol{w}_j}{\|\boldsymbol{w}_j\|} \right\rangle\right)^{p-1} \left(1 - c_{kj}^2 - \left\langle \boldsymbol{\varepsilon}_i, \frac{\boldsymbol{w}_j}{\|\boldsymbol{w}_j\|} \right\rangle c_{kj}\right)$$

$$+\underbrace{\frac{1}{N}\sum_{i\in\mathcal{I}_{k}}y_{i}\;p\left(\left\langle\boldsymbol{\mu}_{k}+\boldsymbol{\varepsilon}_{i},\frac{\boldsymbol{w}_{j}}{\|\boldsymbol{w}_{j}\|}\right\rangle\right)^{p-1}\left\langle\boldsymbol{\mu}_{k},\boldsymbol{\varepsilon}_{i}\right\rangle}_{:=\Gamma_{2}(\text{will be treated later})}$$

$$= \frac{1}{N} \sum_{i \in \mathcal{I}_k} p\left(c_{kj} + \left\langle \boldsymbol{\varepsilon}_i, \frac{\boldsymbol{w}_j}{\|\boldsymbol{w}_j\|} \right\rangle \right)^{p-1} \left(1 - c_{kj}^2 - \left\langle \boldsymbol{\varepsilon}_i, \frac{\boldsymbol{w}_j}{\|\boldsymbol{w}_j\|} \right\rangle c_{kj} \right) + \Gamma_2 \qquad (G.12)$$

With the Taylor expansion

$$\left(c_{kj} + \left\langle \boldsymbol{\varepsilon}_{i}, \frac{\boldsymbol{w}_{j}}{\|\boldsymbol{w}_{j}\|} \right\rangle \right)^{p-1} = c_{kj}^{p-1} + (p-1)c_{kj}^{p-2} \left\langle \boldsymbol{\varepsilon}_{i}, \frac{\boldsymbol{w}_{j}}{\|\boldsymbol{w}_{j}\|} \right\rangle + R_{L} \left| \left\langle \boldsymbol{\varepsilon}_{i}, \frac{\boldsymbol{w}_{j}}{\|\boldsymbol{w}_{j}\|} \right\rangle \right|^{2}, \quad (G.13)$$

where $R_L = \frac{(p-1)(p-2)(c_{kj}+\zeta_L)^{p-2}}{2}$ and ζ_L between 0 and $\left\langle \varepsilon_i, \frac{w_j}{\|w_j\|} \right\rangle$ comes from the Lagrange residual. Clearly $|R_L| \leq 2^{p-3}p^2$. Combining (G.12)(G.13), we have

$$(a) = (G.12)$$

$$= pc_{kj}^{p-1}(1-c_{kj}^{2}) + \left(-pc_{kj}^{p}+p(p-1)c_{kj}^{p-2}(1-c_{kj}^{2})\right)\sum_{i\in\mathcal{I}_{k}}\left\langle\varepsilon_{i},\frac{\boldsymbol{w}_{j}}{\|\boldsymbol{w}_{j}\|}\right\rangle$$
$$+ \frac{1}{N}\left(-p(p-1)c_{kj}+pR_{L}(1-c_{kj}^{2})\right)\sum_{i\in\mathcal{I}_{k}}\left|\left\langle\varepsilon_{i},\frac{\boldsymbol{w}_{j}}{\|\boldsymbol{w}_{j}\|}\right\rangle\right|^{2}$$

$$-\frac{1}{N}pc_{kj}R_L\sum_{i\in\mathcal{I}_k}\left|\left\langle\boldsymbol{\varepsilon}_i,\frac{\boldsymbol{w}_j}{\|\boldsymbol{w}_j\|}\right\rangle\right|^2\left\langle\boldsymbol{\varepsilon}_i,\frac{\boldsymbol{w}_j}{\|\boldsymbol{w}_j\|}\right\rangle+\Gamma_2$$

$$2206 \qquad \geq pc_{kj}^{p-1}(1-c_{kj}^{2}) - \frac{1}{N}p^{2} \left\| \sum_{i \in \mathcal{I}_{k}} \varepsilon_{i} \right\| - \frac{1}{N}2^{p-1}p^{2} \sum_{i \in \mathcal{I}_{k}} \|\varepsilon_{i}\|^{2} - \frac{1}{N}2^{p-3}p^{3} \sum_{i \in \mathcal{I}_{k}} \|\varepsilon_{i}\|^{3} + \Gamma_{2}$$

$$2209 \qquad \geq pc_{kj}^{p-1}(1-c_{kj}^{2}) - \frac{1}{N}p^{2} \left\| \sum_{i \in \mathcal{I}_{k}} \varepsilon_{i} \right\| - \frac{1}{N}2^{p-1}p^{2} \sum_{i \in \mathcal{I}_{k}} \|\varepsilon_{i}\|^{2} - \frac{1}{N}2^{p-3}p^{3} \sum_{i \in \mathcal{I}_{k}} \|\varepsilon_{i}\|^{2} \max_{i} \|\varepsilon_{i}\| + \Gamma_{2}$$

$$2211 \qquad \geq pc_{kj}^{p-1}(1-c_{kj}^{2}) - Cp^{2}\sqrt{\log\frac{K}{\delta}} \frac{\alpha}{\sqrt{N}} - 2^{p-1}p^{3}C^{2}\log\frac{K}{\delta}\alpha^{2} - 2^{p-3}p^{3}C^{3}\log\frac{K^{2}N}{\delta}\alpha^{3} + \Gamma_{2}.$$

$$(G.14)$$

We leave the bound as the last one for now and turn to the other term:

$$(b) = -\frac{1}{N} \sum_{l \neq k} \sum_{i \in \mathcal{I}_{l}: \langle \boldsymbol{x}_{i}, \boldsymbol{w}_{j} \rangle > 0} y_{i} p\left(\left\langle \boldsymbol{x}_{i}, \frac{\boldsymbol{w}_{j}}{\|\boldsymbol{w}_{j}\|}\right\rangle\right)^{p-1} \left(\left\langle \boldsymbol{\mu}_{k}, \boldsymbol{x}_{i}\right\rangle - \left\langle \boldsymbol{x}_{i}, \frac{\boldsymbol{w}_{j}}{\|\boldsymbol{w}_{j}\|}\right\rangle c_{kj}\right)$$
$$= \frac{1}{N} \sum_{l \neq k} \sum_{i \in \mathcal{I}_{l}: \langle \boldsymbol{x}_{i}, \boldsymbol{w}_{j} \rangle > 0} y_{i} p\left(c_{lj} + \left\langle \boldsymbol{\varepsilon}_{i}, \frac{\boldsymbol{w}_{j}}{\|\boldsymbol{w}_{j}\|}\right\rangle\right)^{p} c_{kj}$$
$$\underbrace{-\frac{1}{N} \sum_{l \neq k} \sum_{i \in \mathcal{I}_{l}: \langle \boldsymbol{x}_{i}, \boldsymbol{w}_{j} \rangle > 0} y_{i} p\left(c_{lj} + \left\langle \boldsymbol{\varepsilon}_{i}, \frac{\boldsymbol{w}_{j}}{\|\boldsymbol{w}_{j}\|}\right\rangle\right)^{p-1} \left\langle \boldsymbol{\mu}_{k}, \boldsymbol{\varepsilon}_{i} \right\rangle}_{:=\Gamma_{3}(\text{will be treated later})}$$
$$\geq -\frac{1}{N} \sum_{l \neq k} \sum_{i \in \mathcal{I}_{l}} p\left(|c_{lj}| + \left|\left\langle \boldsymbol{\varepsilon}_{i}, \frac{\boldsymbol{w}_{j}}{\|\boldsymbol{w}_{j}\|}\right\rangle\right|\right)^{p} c_{kj} + \Gamma_{3}$$
(G.15)

With the Taylor expansion

$$\begin{aligned} & \text{Find the traject equation} \\ & \left(|c_{lj}| + \left| \left\langle \epsilon_{i}, \frac{|\mathbf{w}_{j}|}{|\mathbf{w}_{j}|} \right\rangle \right| \right)^{p} = |c_{lj}|^{p} + p|\alpha_{j}|^{p-1} \left| \left\langle \epsilon_{i}, \frac{|\mathbf{w}_{j}|}{|\mathbf{w}_{j}|} \right\rangle \right|^{2} + \mathbb{R}_{L} \left| \left\langle \epsilon_{i}, \frac{|\mathbf{w}_{j}|}{|\mathbf{w}_{j}|} \right\rangle \right|^{2} \\ & \text{where } R_{L} = \frac{p(p-1)(|c_{1}|+\zeta_{L})^{p-2}}{2} \text{ and } \zeta_{L} \text{ between } 0 \text{ and } \left| \left\langle \epsilon_{i}, \frac{|\mathbf{w}_{j}|}{|\mathbf{w}_{j}|} \right\rangle \right| \text{ comes from the Lagrange residual. Clearly } |R_{L}| \leq 2^{p-2}p^{2}. \text{ Combining (G.12)(G.13), we have} \\ & (b) \\ & = (G.15) \end{aligned}$$

$$= -\sum_{l\neq k} p|c_{lj}|^{p}c_{kj} - \frac{1}{N} \sum_{l\neq k} \sum_{i\in \mathbb{Z}_{l}} p(p-1)|c_{lj}|^{p-1}c_{kj} \left| \left\langle \epsilon_{i}, \frac{|\mathbf{w}_{j}|}{||\mathbf{w}_{j}|} \right\rangle \right| - \frac{1}{N} \sum_{l\neq k} \sum_{i\in \mathbb{Z}_{l}} pR_{L}c_{kj} \left| \left\langle \epsilon_{i}, \frac{|\mathbf{w}_{j}|}{||\mathbf{w}_{j}|} \right\rangle \right|^{2} + \Gamma_{3} \end{aligned}$$

$$\geq -\sum_{l\neq k} p|c_{lj}|^{p}c_{kj} - \frac{1}{N} \sum_{l\neq k} \sum_{i\in \mathbb{Z}_{l}} p(p-1)|c_{lj}|^{p-1}c_{kj} \left| \left\langle \epsilon_{i}, \frac{|\mathbf{w}_{j}|}{||\mathbf{w}_{j}|} \right\rangle \right| - K2^{p-2}p^{3}C^{2}\log\frac{K}{\delta}\alpha^{2} + \Gamma_{3}, \end{aligned}$$

$$\geq -\sum_{l\neq k} p|c_{lj}|^{p}c_{kj} - \frac{1}{N} \sum_{l\neq k} \sum_{i\in \mathbb{Z}_{l}} p(p-1)|c_{lj}|^{p-1}c_{kj} \left(||\epsilon_{i}||\sqrt{1-c_{kj}^{2}} + |\langle\epsilon_{i}, \mu_{k}\rangle| \right) - K2^{p-2}p^{3}C^{2}\log\frac{K}{\delta}\alpha^{2} + \Gamma_{3}, \end{aligned}$$

$$\geq -\sum_{l\neq k} p|c_{lj}|^{p}c_{kj} - \sum_{l\neq k} p^{2}|c_{lj}|^{p-1}c_{kj}C\sqrt{\log\frac{K^{2}N}{\delta}}\alpha\sqrt{1-c_{kj}^{2}} - \frac{1}{N}\sum_{l\neq k} \sum_{i\in \mathbb{Z}_{l}} p(p-1)|c_{lj}|^{p-1}c_{kj}\left(||\epsilon_{i}||\sqrt{1-c_{kj}^{2}} + |\langle\epsilon_{i}, \mu_{k}\rangle| \right) - K2^{p-2}p^{3}C^{2}\log\frac{K}{\delta}\alpha^{2} + \Gamma_{3}, \end{aligned}$$

$$\geq -\sum_{l\neq k} p|c_{lj}|^{p}c_{kj} - \sum_{l\neq k} \sum_{i\in \mathbb{Z}_{l}} p(p-1)|c_{lj}|^{p-1}c_{kj}\left(\sqrt{\log\frac{K^{2}N}{\delta}}\alpha\sqrt{1-c_{kj}^{2}} - \sum p^{2}\beta^{2}C\log\frac{K}{\delta}\alpha^{2} + \Gamma_{3} + \frac{1}{N}\sum_{l\neq k} \sum_{i\in \mathbb{Z}_{l}} p(p-1)|c_{lj}|^{p-1}c_{kj}\left(||\epsilon_{i}||\sqrt{1-c_{kj}^{2}} - K2^{p-2}p^{3}C^{2}\log\frac{K}{\delta}\alpha^{2} + \Gamma_{3} + \frac{1}{N}\sum_{l\neq k} \sum_{i\in \mathbb{Z}_{l}} p(p-1)|c_{lj}|^{p-1}c_{kj}\left(\sqrt{\log\frac{K^{2}N}{\delta}}\alpha\sqrt{1-c_{kj}^{2}} \right) - K2^{p-2}p^{3}C^{2}\log\frac{K}{\delta}\alpha^{2} + \Gamma_{3} + \Gamma_{4} + \frac{1}{N}\sum_{l\neq k} \sum_{i\in \mathbb{Z}_{l}} p(p-1)|c_{lj}|^{p-1}c_{kj}\left(\sqrt{\log\frac{K^{2}N}{\delta}}\alpha\sqrt{1-c_{kj}^{2}} \right) - K2^{p-2}p^{3}C^{2}\log\frac{K}{\delta}\alpha^{2} + \Gamma_{3} + \Gamma_{4} + \frac{1}{N}\sum_{l\neq k} \sum_{i\in \mathbb{Z}_{l}} p(p-1)\left(\frac{1}{N}\sum_{l\neq$$

$$\geq -pc_{kj}^{p-1} \left(\max_{l \neq k} \frac{|c_{lj}|^{p-2}}{c_{kj}^{p-2}} (1 - c_{kj}^2) \right) \left(1 - pC\sqrt{K \log \frac{K^2 N}{\delta}} \alpha (1 - c_{kj}^2) \right) - K2^{p-2} p^3 C^2 \log \frac{K}{\delta} \alpha^2 + \Gamma_3 + \Gamma_4$$

$$\geq -\frac{p}{2} c_{kj}^{p-1} \left(\max_{l \neq k} \frac{|c_{lj}|^{p-2}}{c_{kj}^{p-2}} \right) (1 - c_{kj}^2) - K2^{p-2} p^3 C^2 \log \frac{K}{\delta} \alpha^2 + \Gamma_3 + \Gamma_4$$

$$(G.17)$$

Finally, combining (G.14)(G.17), we have

2275
2276

$$\frac{d}{dt}c_{kj} \ge pc_{kj}^{p-1} \left(1 - \max_{l \ne k} \frac{|c_{lj}|^{p-2}}{c_{kj}^{p-2}}\right) (1 - c_{kj}^2) - C_1' \sqrt{\log \frac{K}{\delta}} \frac{\alpha}{\sqrt{N}} - C_2' \log \frac{K}{\delta} \alpha^2 - C_3' \log \frac{K^2 N}{\delta} \alpha^3$$
2278

$$- |\Gamma_1| - |\Gamma_2| - |\Gamma_3| - |\Gamma_4|$$
2279

$$(1 - c_{kj}^2) - C_1' \sqrt{\log \frac{K}{\delta}} \frac{\alpha}{\sqrt{N}} - C_2' \log \frac{K}{\delta} \alpha^2 - C_3' \log \frac{K^2 N}{\delta} \alpha^3$$

$$\geq pc_{kj}^{p-1} \left(1 - \max_{l \neq k} \frac{|c_{lj}|^{p-2}}{c_{kj}^{p-2}} \right) (1 - c_{kj}) - C_1' \sqrt{\log \frac{K}{\delta}} \frac{\alpha}{\sqrt{N}} - C_2' \log \frac{K}{\delta} \alpha^2 - C_3' \log \frac{K^2 N}{\delta} \alpha^3 - |\Gamma_1| - |\Gamma_2| - |\Gamma_3| - |\Gamma_4|,$$

where the readers should be able to find universal constants C'_1, C'_2, C'_3 from the derivation. It remains to bound these $|\Gamma_i|, i = 1, \dots, 4$. Indeed, we can find the following bound:

$$\begin{split} |\Gamma_{1}| &= \left| \frac{1}{N} \sum_{i:\langle \boldsymbol{x}_{i}, \boldsymbol{w}_{j} \rangle > 0} f^{(p)}(\boldsymbol{x}_{i}; \boldsymbol{\theta}) p\left(\left\langle \boldsymbol{x}_{i}, \frac{\boldsymbol{w}_{j}}{\|\boldsymbol{w}_{j}\|} \right\rangle \right)^{p-1} \left(\left\langle \boldsymbol{\mu}_{k}, \boldsymbol{x}_{i} \right\rangle - \left\langle \boldsymbol{x}_{i}, \frac{\boldsymbol{w}_{j}}{\|\boldsymbol{w}_{j}\|} \right\rangle c_{kj} \right) \right| \\ &\leq \max_{i} |f^{(p)}(\boldsymbol{x}_{i}; \boldsymbol{\theta})| \frac{1}{N} \sum_{i:\langle \boldsymbol{x}_{i}, \boldsymbol{w}_{j} \rangle > 0} |p| \|\boldsymbol{x}_{i}\|^{p-1} \left(2 \|\boldsymbol{x}_{i}\| \right)| \\ &\leq p 2^{p+1} \max_{i} |f^{(p)}(\boldsymbol{x}_{i}; \boldsymbol{\theta})|, \\ |\Gamma_{2}| &= \left| \frac{1}{N} \sum_{i \in \mathcal{I}_{k}} y_{i} p\left(\left\langle \boldsymbol{\mu}_{k} + \boldsymbol{\varepsilon}_{i}, \frac{\boldsymbol{w}_{j}}{\|\boldsymbol{w}_{j}\|} \right\rangle \right)^{p-1} \left\langle \boldsymbol{\mu}_{k}, \boldsymbol{\varepsilon}_{i} \right\rangle \right| \\ &\leq \frac{1}{N} \sum_{i \in \mathcal{I}_{k}} p \|\boldsymbol{x}_{i}\|^{p-1} |\langle \boldsymbol{\mu}_{k}, \boldsymbol{\varepsilon}_{i} \rangle| \leq p 2^{p-1} C \sqrt{\log \frac{K^{2}N}{\delta}} \frac{\alpha}{\sqrt{D}} \\ |\Gamma_{3}| &= \left| -\frac{1}{N} \sum_{l \neq k} \sum_{i \in \mathcal{I}_{l}: \langle \boldsymbol{x}_{i}, \boldsymbol{w}_{j} \rangle > 0} y_{i} p\left(c_{lj} + \left\langle \boldsymbol{\varepsilon}_{i}, \frac{\boldsymbol{w}_{j}}{\|\boldsymbol{w}_{j}\|} \right\rangle \right)^{p-1} \left\langle \boldsymbol{\mu}_{k}, \boldsymbol{\varepsilon}_{i} \right\rangle \right| \\ &\leq \frac{1}{N} \sum_{i \in \mathcal{I}_{k}} p \|\boldsymbol{x}_{i}\|^{p-1} |\langle \boldsymbol{\mu}_{k}, \boldsymbol{\varepsilon}_{i} \rangle| \leq p 2^{p-1} C \sqrt{\log \frac{K^{2}N}{\delta}} \frac{\alpha}{\sqrt{D}} \\ |\Gamma_{4}| &= \left| -\frac{1}{N} \sum_{l \neq k} \sum_{i \in \mathcal{I}_{l}} p(p-1) |c_{lj}|^{p-1} c_{kj} |\langle \boldsymbol{\varepsilon}_{i}, \boldsymbol{\mu}_{k} \rangle| \right| \\ &\leq \frac{1}{N} \sum_{l \neq k} \sum_{i \in \mathcal{I}_{k}} p^{2} |\langle \boldsymbol{\mu}_{k}, \boldsymbol{\varepsilon}_{i} \rangle| \leq K p^{2} C \sqrt{\log \frac{K^{2}N}{\delta}} \frac{\alpha}{\sqrt{D}}. \end{split}$$

With these norm bounds, we have

$$\begin{split} \frac{d}{dt} c_{kj} \\ \geq p c_{kj}^{p-1} \left(1 - \max_{l \neq k} \frac{|c_{lj}|^{p-2}}{c_{kj}^{p-2}} \right) (1 - c_{kj}^2) - C_1' \sqrt{\log \frac{K}{\delta}} \frac{\alpha}{\sqrt{N}} - C_2' \log \frac{K}{\delta} \alpha^2 \\ - C_3' \log \frac{K^2 N}{\delta} \alpha^3 - C_4' \max_i |f^{(p)}(\boldsymbol{x}_i; \boldsymbol{\theta})| - C_5' \sqrt{\log \frac{K^2 N}{\delta}} \frac{\alpha}{\sqrt{D}} \end{split}$$

$$\geq p c_{kj}^{p-1} \left(1 - \max_{l \neq k} \frac{|c_{lj}|^{p-2}}{c_{kj}^{p-2}} \right) (1 - c_{kj}) - C_1 \log \frac{K}{\delta} \alpha^2 - C_2 \max_i |f^{(p)}(\boldsymbol{x}_i; \boldsymbol{\theta})|.$$

Lastly, our bound is

1. When we only assumed $c_{kj} \geq \tilde{\Delta}_1$:

$$\frac{d}{dt}c_{kj} \ge pc_{kj}^{p-1} \left(1 - \max_{l \neq k} \frac{|c_{lj}|^{p-2}}{c_{kj}^{p-2}}\right) (1 - c_{kj}) - C_1 \log \frac{K}{\delta} \alpha^2 - C_2 \max_i |f^{(p)}(\boldsymbol{x}_i; \boldsymbol{\theta})|$$
$$\ge pc_{kj}^{p-1} \tilde{\Delta}_2 (1 - c_{kj}) - C_1 \log \frac{K}{\delta} \alpha^2 - C_2 \max_i |f^{(p)}(\boldsymbol{x}_i; \boldsymbol{\theta})|$$

2. When we further assume $c_{kj} \ge \sqrt{\frac{4}{5}}$, we have that $\sum_{l \ne k} c_{lj}^2 = 1 - c_{kj}^2 \le \frac{1}{5}$, then $\max_{l \ne k} |c_{lj}| \le \sqrt{\frac{1}{5}}$. Therefore

$$\left(1 - \max_{l \neq k} \frac{|c_{lj}|^{p-2}}{c_{kj}^{p-2}}\right) = \left(1 - \left(\frac{\max_{l \neq k} |c_{lj}|}{c_{kj}}\right)^{p-2}\right) \ge 1 - \frac{1}{2^{p-2}}, \quad (G.18)$$

which leads to

$$\frac{d}{dt}c_{kj} \ge pc_{kj}^{p-1} \left(1 - \frac{1}{2^{p-2}}\right) (1 - c_{kj}) - C_1 \log \frac{K}{\delta} \alpha^2 - C_2 \max_i |f^{(p)}(\boldsymbol{x}_i; \boldsymbol{\theta})|.$$

Lemma 7 (Restated). Let p > 2. Condition on good event \mathcal{E}_{good} . Given an initialization shape that satisfies Assumption 2 with non-degeneracy gap $\Delta > 0$, define

$$t_{1a} := \inf \left\{ t : \max_{i} |f^{(p)}(\boldsymbol{x}_{i}; \boldsymbol{\theta}(t)| > \min \left\{ \frac{\tilde{\Delta}_{1}^{p-1} \tilde{\Delta}_{2}(1-\tilde{\Delta}_{1})}{2^{p+1}}, \frac{\tilde{\Delta}_{1}^{p-1} \tilde{\Delta}_{2}(1-\sqrt{2\Delta})}{2K2^{p+1}} \right\} \right\}.$$
(G.19)

2351 Then the following holds $\forall t \leq t_{1a}$:

$$c_{kj}(t) \ge c_{kj}(0) \ge \tilde{\Delta}_1, \forall 1 \le k \le K, j \in \mathcal{N}_k , \qquad (G.20)$$

and

$$\frac{|c_{lj}^{p-2}(t)|}{c_{kj}^{p-2}(t)} \le \frac{|c_{lj}^{p-2}(0)|}{c_{kj}^{p-2}(0)} \le 1 - \tilde{\Delta}_2 \text{ and } \forall l \neq k, j \in \mathcal{N}_k \text{.}$$
(G.21)

Proof. When $1 \le k \le K_1$, $j \in \mathcal{N}_k$ implies that $j \in \mathcal{N}_+$ thus $\operatorname{sign}(v_j) = 1$. We shall primarily focus on this case as the proof is nearly identical for $K_1 + 1 \le k \le K$.

Overview of the proof: We will prove by contradiction, we let $\tau_1 := \inf\{t : \exists k, j \in \mathcal{N}_k, s.t. c_{kj}(t) < c_{kj}(0)\}$ and $\tau_2 := \inf\{t : \exists k, j \in \mathcal{N}_k, \&l \neq k, s.t. \frac{|c_{lj}^{p-2}(t)|}{c_{kj}(t)} > \frac{|c_{lj}^{p-2}(0)|}{c_{kj}(0)}\}$, by the continuity of every $c_{kj}(t)$ and every $\frac{|c_{lj}^{p-2}(t)|}{c_{kj}(t)}$ on the interval $[0, \tau_1]$ and $[0, \tau_2]$ respectively, we know that $c_{kj}(\tau_1) = c_{kj}(0)$ for some k, j and $\frac{|c_{lj}^{p-2}(\tau_2)|}{c_{kj}(\tau_2)} = \frac{|c_{lj}^{p-2}(0)|}{c_{kj}(0)}$ for some k, j, l. If $\min\{\tau_1, \tau_2\} > t_{1a}$ then there is nothing to be proved, otherwise, there are two cases:

1. When
$$\tau_1 = \min\{\tau_1, \tau_2\} \le t_{1a}$$
, we show that for the k, j such that $c_{kj}(\tau_1) = c_{kj}(0)$
$$\frac{d}{dt}c_{kj}\Big|_{t=-} \ge 0, \qquad (G.22)$$

which says $c_{kj}(\tau_1 + \Delta t) \ge c_{kj}(0)$ for every sufficiently small Δt , contradicting the definition of τ_1 .

2377 2. When $\tau_2 = \min\{\tau_1, \tau_2\} \le t_{1a}$, we show that for the k, j, l such that $\frac{|c_{lj}^{p-2}(\tau_2)|}{c_{kj}(\tau_2)} = \frac{|c_{lj}^{p-2}(0)|}{c_{kj}(0)}$ 2379 $\frac{d}{dt} \log \frac{|c_{lj}|}{c_{kj}}\Big|_{t=\tau_2} \le 0$, (G.23)

which says $\frac{|c_{lj}(\tau_1+\Delta t)|}{c_{kj}(\tau_1+\Delta t)} \leq c_{kj}(0)$ for every sufficiently small Δt (due to the monotonicity of log function), contradicting the definition of τ_2 .

Time derivatives of log cosine angles We have shown in (D.6) that for every $1 \le l \le D$, whenever $|c_{lj}| > 0$,

$$\frac{d}{dt} \log |c_{lj}|$$

$$= -\frac{1}{N} \sum_{i: \langle \boldsymbol{x}_i, \boldsymbol{w}_j \rangle > 0} \nabla_{\hat{y}} \ell_i \, p\left(\left\langle \boldsymbol{x}_i, \frac{\boldsymbol{w}_j}{\|\boldsymbol{w}_j\|} \right\rangle\right)^{p-1} \left(\frac{\langle \boldsymbol{\mu}_l, \boldsymbol{x}_i \rangle}{c_{lj}} - \left\langle \boldsymbol{x}_i, \frac{\boldsymbol{w}_j}{\|\boldsymbol{w}_j\|} \right\rangle\right)$$

$$= -\frac{1}{N} \sum_{i:\langle \boldsymbol{x}_i, \boldsymbol{w}_j \rangle > 0} \nabla_{\hat{y}} \ell_i \ p\left(\left\langle \boldsymbol{x}_i, \frac{\boldsymbol{w}_j}{\|\boldsymbol{w}_j\|} \right\rangle\right)^{p-1} \frac{\langle \boldsymbol{\mu}_l, \boldsymbol{x}_i \rangle}{c_{lj}} \\ + \frac{1}{N} \sum_{i:\langle \boldsymbol{x}_i, \boldsymbol{w}_j \rangle > 0} \nabla_{\hat{y}} \ell_i \ p\left(\left\langle \left\langle \boldsymbol{x}_i, \frac{\boldsymbol{w}_j}{\|\boldsymbol{w}_j\|} \right\rangle\right)^p \ .$$

Case One: $\tau_1 = \min{\{\tau_1, \tau_2\}}$. This case is relatively easier as we have already shown Lemma 6. For the k, j such that $c_{kj} = \tilde{\Delta}$

$$\frac{d}{dt}c_{kj}\Big|_{t=\tau_1} \ge pc_{kj}^{p-1}\tilde{\Delta}_2(1-c_{kj}) - C_1\log\frac{K}{\delta}\alpha^2 - \underbrace{p2^{p+1}\max_i|f^{(p)}(\boldsymbol{x}_i;\boldsymbol{\theta})|}_{(s)}$$

by Lemma 6 (conditions are satisified at $t = \tau_1$ and one should be able to get (*) using the intermediate results in the proof of Lemma 6). Then

$$\frac{d}{dt}c_{kj}\Big|_{t=\tau_1} \ge p\tilde{\Delta}_1^{p-1}\tilde{\Delta}_2(1-\tilde{\Delta}_1) - C_1\log\frac{K}{\delta}\alpha^2 - p2^{p+1}\max_i |f^{(p)}(\boldsymbol{x}_i;\boldsymbol{\theta})|,$$

$$\stackrel{(\tau_1 \le t_{1a})}{\ge} p\tilde{\Delta}_1^{p-1}\tilde{\Delta}_2(1-\tilde{\Delta}_1) - C_1\log\frac{K}{\delta}\alpha^2 - \frac{1}{2}p\tilde{\Delta}_1^{p-1}\tilde{\Delta}_2(1-\tilde{\Delta}_1)$$

$$\ge \frac{1}{2}p\tilde{\Delta}_1^{p-1}\tilde{\Delta}_2(1-\tilde{\Delta}_1) - C_1\log\frac{K}{\delta}\alpha^2 \ge 0,$$

for sufficiently small α .

Case Two: $\tau_2 = \min{\{\tau_1, \tau_2\}}$. For the k, j, l such that $\frac{|c_{lj}^{p-2}(\tau_2)|}{c_{kj}(\tau_2)} = \frac{|c_{lj}^{p-2}(0)|}{c_{kj}(0)}$, we have (although we omit the notation, all the derivations are at τ_2 , so that c_{lj} can appear in the denominator of a fraction.)

$$\begin{aligned} \frac{d}{dt} \log \frac{|c_{lj}|}{c_{kj}} \\ \frac{2419}{2419} & \frac{d}{dt} \log \frac{|c_{lj}|}{c_{kj}} \\ \frac{2420}{2421} & = \frac{d}{dt} \log |c_{lj}| - \frac{d}{dt} \log c_{kj} \\ \frac{2422}{2423} & = -\frac{1}{N} \sum_{i: \langle \boldsymbol{x}_i, \boldsymbol{w}_j \rangle > 0} \nabla_{j} \ell_i \, p \left(\left\langle \boldsymbol{x}_i, \frac{\boldsymbol{w}_j}{\|\boldsymbol{w}_j\|} \right\rangle \right)^{p-1} \left(\frac{\langle \boldsymbol{\mu}_l, \boldsymbol{x}_i \rangle}{c_{lj}} - \frac{\langle \boldsymbol{\mu}_k, \boldsymbol{x}_i \rangle}{c_{kj}} \right) \\ \frac{2425}{2426} & = \frac{1}{N} \sum_{i: \langle \boldsymbol{x}_i, \boldsymbol{w}_j \rangle > 0} (y_i - f^{(p)}(\boldsymbol{x}_i; \boldsymbol{\theta})) \, p \left(\left\langle \boldsymbol{x}_i, \frac{\boldsymbol{w}_j}{\|\boldsymbol{w}_j\|} \right\rangle \right)^{p-1} \left(\frac{\langle \boldsymbol{\mu}_l, \boldsymbol{x}_i \rangle}{c_{lj}} - \frac{\langle \boldsymbol{\mu}_k, \boldsymbol{x}_i \rangle}{c_{kj}} \right) \\ \frac{2428}{2429} & = \frac{1}{N} \sum_{i: \langle \boldsymbol{x}_i, \boldsymbol{w}_j \rangle > 0} y_i \, p \left(\left\langle \boldsymbol{x}_i, \frac{\boldsymbol{w}_j}{\|\boldsymbol{w}_j\|} \right\rangle \right)^{p-1} \left(\frac{\langle \boldsymbol{\mu}_l, \boldsymbol{x}_i \rangle}{c_{lj}} - \frac{\langle \boldsymbol{\mu}_k, \boldsymbol{x}_i \rangle}{c_{kj}} \right) \end{aligned}$$

 We view Γ_1, Γ_2 as "perturbation term" and will control their norms later. For the first two terms in (G.24), we have, respectively:

$$\frac{1}{N} \sum_{i \in \mathcal{I}_k: \langle \boldsymbol{x}_i, \boldsymbol{w}_j \rangle > 0} y_i p \left(c_{kj} + \left\langle \boldsymbol{\varepsilon}_i, \frac{\boldsymbol{w}_j}{\|\boldsymbol{w}_j\|} \right\rangle \right)^{p-1} \left(-\frac{1}{c_{kj}} \right)$$

$$= -\frac{p}{N} \sum_{i \in \mathcal{I}_k} \left(c_{kj} + \left\langle \boldsymbol{\varepsilon}_i, \frac{\boldsymbol{w}_j}{\|\boldsymbol{w}_j\|} \right\rangle \right)^{p-1} \frac{1}{c_{kj}}$$

$$= -\frac{p}{N} \sum_{i \in \mathcal{I}_k} c_{kj}^{p-2} \left(1 - \frac{\left| \left\langle \boldsymbol{\varepsilon}_i, \frac{\boldsymbol{w}_j}{\|\boldsymbol{w}_j\|} \right\rangle \right|}{c_{kj}} \right)^{p-1}$$

$$\leq -\frac{p}{N} \sum_{i \in \mathcal{I}_k} c_{kj}^{p-2} \left(1 - \frac{\left| \left\langle \boldsymbol{\varepsilon}_i, \frac{\boldsymbol{w}_j}{\|\boldsymbol{w}_j\|} \right\rangle \right|}{c_{kj}} \right)^{p-1} \\ \leq -\frac{p}{N} \sum_{i \in \mathcal{I}_k} c_{kj}^{p-2} \left(1 - (p-1) \frac{\left| \left\langle \boldsymbol{\varepsilon}_i, \frac{\boldsymbol{w}_j}{\|\boldsymbol{w}_j\|} \right\rangle \right|}{c_{kj}} \right)^{p-1} \right)^{p-1}$$

$$\leq -\frac{p}{N} \sum_{i \in \mathcal{I}_k} c_{kj}^{p-2} \left(1 - (p-1) \frac{\left| \left\langle \varepsilon \right\rangle \right|}{2} \right)$$

$$\leq -pc_{kj}^{p-2} + p(p-1)\frac{\max_i \|\boldsymbol{\varepsilon}_i\|}{c_{kj}(0)} \geq -pc_{kj}^{p-2} + p(p-1)\sqrt{8\log\frac{4K^2N}{\delta}}\alpha\,,$$

 c_{kj}

and similarly,

$$\frac{1}{N} \sum_{i \in \mathcal{I}_l: \langle \boldsymbol{x}_i, \boldsymbol{w}_j \rangle > 0} y_i \, p \left(c_{lj} + \left\langle \boldsymbol{\varepsilon}_i, \frac{\boldsymbol{w}_j}{\|\boldsymbol{w}_j\|} \right\rangle \right)^{p-1} \left(\frac{1}{c_{lj}} \right)^{p-1} \left(\frac{1}{$$

Therefore we have

$$\frac{d}{dt}\log\frac{|c_{lj}|}{c_{kj}} \le -p(c_{kj}^{p-2} - |c_{lj}|^{p-2}\mathbb{1}_{l\le K}) + 2p(p-1)\sqrt{8\log\frac{4K^2N}{\delta}}\alpha - |\Gamma_1| - |\Gamma_2|$$
$$\le -p(c_{kj}^{p-2} - |c_{lj}|^{p-2}) + 2p(p-1)\sqrt{8\log\frac{4K^2N}{\delta}}\alpha - |\Gamma_1| - |\Gamma_2|$$
$$\le -pc_{kj}^{p-2}\left(1 - \frac{|c_{lj}|^{p-2}}{c_{kj}^{p-2}}\right) + 2p(p-1)\sqrt{8\log\frac{4K^2N}{\delta}}\alpha - |\Gamma_1| - |\Gamma_2|$$

 $p\!-\!1$

$$\leq -p\tilde{\Delta}_1^{p-2}\tilde{\Delta}_2 + 2p(p-1)\sqrt{8\log\frac{4K^2N}{\delta}\alpha} - |\Gamma_1| - |\Gamma_2|.$$

It remains to bound these $|\Gamma_1|, |\Gamma_2|$. Indeed, we can find the following bound⁴ (note that at τ_2 , we have $|c_{lj}| = c_{kj}(1 - \sqrt{2\Delta}))$ and $c_{kj} \ge \tilde{\Delta}_1$):

$$\begin{split} |\Gamma_1| &= \left| \frac{1}{N} \sum_{i: \langle \boldsymbol{x}_i, \boldsymbol{w}_j \rangle > 0} f^{(p)}(\boldsymbol{x}_i; \boldsymbol{\theta}) \ p\left(\left\langle \boldsymbol{x}_i, \frac{\boldsymbol{w}_j}{\|\boldsymbol{w}_j\|} \right\rangle \right)^{p-1} \left(\frac{\langle \boldsymbol{\mu}_l, \boldsymbol{x}_i \rangle}{c_{lj}} - \frac{\langle \boldsymbol{\mu}_k, \boldsymbol{x}_i \rangle}{c_{kj}} \right) \right| \\ &\leq \left| \frac{1}{N} \sum_{1 \leq i \leq KN} |f^{(p)}(\boldsymbol{x}_i; \boldsymbol{\theta})| \ p \left| \left\langle \boldsymbol{x}_i, \frac{\boldsymbol{w}_j}{\|\boldsymbol{w}_j\|} \right\rangle \right|^{p-1} \left(\frac{|\langle \boldsymbol{\mu}_l, \boldsymbol{x}_i \rangle|}{c_{kj}(1 - \sqrt{2\Delta})} + \frac{|\langle \boldsymbol{\mu}_k, \boldsymbol{x}_i \rangle|}{c_{kj}} \right) \right| \end{split}$$

$$\leq \max_{i} |f^{(p)}(\boldsymbol{x}_{i};\boldsymbol{\theta})| \frac{Kp2^{p+1}}{\tilde{\Delta}_{1}(1-\sqrt{2\Delta})} \stackrel{(\tau_{2}\leq t_{1a})}{\leq} \frac{1}{2}p\tilde{\Delta}_{1}^{p-2}\tilde{\Delta}_{2},$$

$$|\Gamma_2| = \left| \frac{1}{N} \sum_{i: \langle \boldsymbol{x}_i, \boldsymbol{w}_j \rangle > 0} y_i \, p\left(\left\langle \boldsymbol{x}_i, \frac{\boldsymbol{w}_j}{\|\boldsymbol{w}_j\|} \right\rangle \right)^{p-1} \left(\frac{\langle \boldsymbol{\mu}_l, \boldsymbol{\varepsilon}_i \rangle}{c_{lj}} - \frac{\langle \boldsymbol{\mu}_k, \boldsymbol{\varepsilon}_i \rangle}{c_{kj}} \right) \right|$$

$$\leq \left| \frac{1}{N} \sum_{1 \leq i \leq KN} p \left| \left\langle \boldsymbol{x}_{i}, \frac{\boldsymbol{w}_{j}}{\|\boldsymbol{w}_{j}\|} \right\rangle \right|^{p-1} \left(\frac{|\langle \boldsymbol{\mu}_{l}, \boldsymbol{\varepsilon}_{i} \rangle|}{c_{kj}(1 - \sqrt{2\Delta})} + \frac{|\langle \boldsymbol{\mu}_{k}, \boldsymbol{\varepsilon}_{i} \rangle|}{c_{kj}} \right) \right. \\ \leq \max_{i,k} \left| \left\langle \boldsymbol{\mu}_{k}, \boldsymbol{\varepsilon}_{i} \right\rangle \left| \frac{Kp2^{p+1}}{\tilde{\Delta}_{1}(1 - \sqrt{2\Delta})} \leq \frac{CKp2^{p+1}}{\tilde{\Delta}_{1}(1 - \sqrt{2\Delta})} \sqrt{\log \frac{K^{2}N}{\delta}} \frac{\alpha}{\sqrt{D}} \right.$$

Finally, we arrived at

$$\frac{d}{dt} \log \frac{|c_{lj}|}{c_{kj}}\Big|_{t=\tau_2}$$

$$\leq -p\tilde{\Delta}_1^{p-2}\tilde{\Delta}_2 + 2p(p-1)\sqrt{8\log\frac{4K^2N}{\delta}\alpha} + \frac{1}{2}p\tilde{\Delta}_1^{p-2}\tilde{\Delta}_2 + \frac{CKp2^{p+1}}{\tilde{\Delta}_1(1-\sqrt{2\Delta})}\sqrt{\log\frac{K^2N}{\delta}\frac{\alpha}{\sqrt{D}}}$$

⁴It may take some time to recollect the terms we omitted in (G.24) and regroup them into Γ_2

2538
2539
$$\leq -\frac{1}{2}p\tilde{\Delta}_{1}^{p-2}\tilde{\Delta}_{2} + 2p(p-1)\sqrt{8\log\frac{4K^{2}N}{\delta}}\alpha + \frac{CKp2^{p+1}}{\tilde{\Delta}_{1}(1-\sqrt{2\Delta})}\sqrt{\log\frac{K^{2}N}{\delta}}\frac{\alpha}{\sqrt{D}} \leq 0, \quad (G.25)$$
2540

for sufficiently small α .

Lemma 8 (Restated). Let p > 2. Condition on good event \mathcal{E}_{good} , then with any balanced initialization scale $\epsilon \leq \frac{1}{4\sqrt{h}W_{\max}^2}$, the solution to gradient flow dynamics satisfies

$$\max_{k} |f^{(p)}(\boldsymbol{\mu}_{k};\boldsymbol{\theta}(t))| \leq 2\epsilon \sqrt{h} W_{\max}^{2}, \quad \forall t \leq \frac{1}{2^{p+2}K} \log\left(\frac{1}{2^{p-1}\sqrt{h\epsilon}}\right).$$
(G.26)

Proof. Let $T := \inf\{t : \max_i | f(\boldsymbol{x}_k; \boldsymbol{\theta}(t))| > 2\epsilon \sqrt{h} W_{\max}^2\}$, then $\forall t \leq T, j \in [h]$, we have

$$\frac{d}{dt} \|\boldsymbol{w}_j\|^2 = -2 \frac{\operatorname{sign}(v_j(0))}{N} \left(\sum_{i: \langle \boldsymbol{x}_i, \boldsymbol{w}_j \rangle > 0} \nabla_{\hat{y}} \ell_i \left(\left\langle \boldsymbol{x}_i, \frac{\boldsymbol{w}_j}{\|\boldsymbol{w}_j\|} \right\rangle \right)^p \right) \|\boldsymbol{w}_j\|^2$$

$$\leq 2rac{1}{N}\sum_{i=1}^{KN}|
abla_{\hat{y}}\ell_i|\|oldsymbol{w}_j\|^2rac{(\langleoldsymbol{x}_i,oldsymbol{w}_j
angle^p}{\|oldsymbol{w}_j\|^p}$$

$$\leq 2rac{1}{N}\sum_{i=1}^{KN}|
abla_{\hat{y}}\ell_i|\|oldsymbol{w}_j\|^2\|oldsymbol{x}_i\|^p \ 2^{n+1}|KN|$$

$$\leq rac{2^{p+1}}{N} \sum_{i=1}^{KN} (1 + |f(m{x}_k;m{ heta}(t))|) \|m{w}_j\|^2$$

2562
2563
2564
2565

$$\leq \frac{2^{p+1}}{N} \sum_{i=1}^{KN} (1 + 4\epsilon \sqrt{h} W_{\max}^2) \|\boldsymbol{w}_j\|^2$$
2565

$$\leq 2^{p+1} V(1 + 4\epsilon \sqrt{h} W_{\max}^2) \|\boldsymbol{w}_j\|^2$$

$$\leq 2^{p+1} K (1 + 4\epsilon \sqrt{h} W_{\max}^2) \|\boldsymbol{w}_j\|^2.$$
(G.27)

Let $\tau_j := \inf\{t : \|\boldsymbol{w}_j(t)\|^2 > \frac{2\epsilon M^2}{2^{p-1}\sqrt{h}}\}$, and let $j^* := \arg\min_j \tau_j$, then $\tau_{j^*} = \min_j \tau_j \leq T$ due to the fact that

$$|f(\boldsymbol{x}_i;\boldsymbol{\theta})| = \left|\sum_{j\in[h]} \mathbb{1}_{\langle \boldsymbol{w}_j, \boldsymbol{x}_i \rangle > 0} v_j \frac{(\langle \boldsymbol{w}_j, \boldsymbol{x}_k \rangle)^p}{\|\boldsymbol{w}_j\|^{p-1}}\right| \le 2^p \sum_{j\in[h]} \|\boldsymbol{w}_j\|^2 \le 2^p h \max_{j\in[h]} \|\boldsymbol{w}_j\|^2 \,,$$

which implies $||f(\boldsymbol{x}_k; \boldsymbol{\theta}(t))| > 2\epsilon \sqrt{h} W_{\max}^2 \Rightarrow \exists j, s.t. \| \boldsymbol{w}_j(t) \|^2 > \frac{\epsilon W_{\max}^2}{2^{p-1}\sqrt{h}}$ ".

Then for $t \leq \tau_{j^*}$, we have

$$\frac{d}{dt} \|\boldsymbol{w}_{j^*}\|^2 \le 2^{p+1} K (+4\epsilon \sqrt{h} W_{\max}^2) \|\boldsymbol{w}_{j^*}\|^2.$$
(G.28)

By Grönwall's inequality, we have $\forall t \leq \tau_{j^*}$

2580
2581
$$\|\boldsymbol{w}_{j^*}(t)\|^2 \leq \exp\left(2^{p+1}K(1+4\epsilon\sqrt{h}W_{\max}^2)t\right)\|\boldsymbol{w}_{j^*}(0)\|^2,$$
2582
2583
$$= \exp\left(2^{p+1}K(1+4\epsilon\sqrt{h}W_{\max}^2)t\right)\epsilon^2\|\boldsymbol{w}_{j^*0}\|^2$$
2584
2585
$$\leq \exp\left(2^{p+1}K(1+4\epsilon\sqrt{h}W_{\max}^2)t\right)\epsilon^2W_{\max}^2.$$

Suppose $\tau_{j^*} < \frac{1}{2^{p+2}K} \log\left(\frac{1}{2^{p-1}\sqrt{h\epsilon}}\right)$, then by the continuity of $||w_{j^*}(t)||^2$, we have

2588
2589
2589
2590
2591

$$\frac{2\epsilon W_{\max}^2}{2^{p-1}\sqrt{h}} \le \|\boldsymbol{w}_{j^*}(\tau_{j^*})\|^2 \le \exp\left(2^{p+1}K(1+4\epsilon\sqrt{h}W_{\max}^2)\tau_{j^*}\right)\epsilon^2 W_{\max}^2$$

$$\le \exp\left(2^{p+1}K(1+4\epsilon\sqrt{h}W_{\max}^2)\frac{1}{2^{p+2}K}\log\left(\frac{1}{2^{p-1}\sqrt{h}\epsilon}\right)\right)\epsilon^2 W_{\max}^2$$

$$2592$$

$$2593$$

$$2594$$

$$2594$$

$$2595$$

$$2596$$

$$\leq \exp\left(\frac{1+4\epsilon\sqrt{h}W_{\max}^2}{2}\log\left(\frac{1}{2^{p-1}\sqrt{h}\epsilon}\right)\right)\epsilon^2 W_{\max}^2$$

$$\leq \exp\left(\log\left(\frac{1}{2^{p-1}\sqrt{h}\epsilon}\right)\right)\epsilon^2 W_{\max}^2 = \frac{\epsilon W_{\max}^2}{2^{p-1}\sqrt{h}},$$

which leads to a contradiction $2\epsilon \leq \epsilon$. Therefore, one must have $T \geq \tau_{j^*} \geq \frac{1}{2^{p+2}K} \log\left(\frac{1}{2^{p-1}\sqrt{h\epsilon}}\right)$. This finishes the proof.

Lemma 9 (Restated). Let p > 2. Given some C > 0, if for some z(t), the following holds

$$\frac{d}{dt}z \ge Cz^{p-1}, \forall t \in [0,T], \ z(0) = z_0, \ z(T) = z_1,$$
(G.29)

for some $0 < z_0 \leq z_1 < 1$. Then the travel time T for z(t) to go from z_0 to z_1 satisfies:

$$T \le \frac{1}{(p-2)Cz_0^{p-2}}.$$
(G.30)

Proof. We have

$$\int_{z_0}^{z_1} \frac{1}{Cz^{p-1}} dz \ge \int_0^T dt \,, \tag{G.31}$$

thus

$$T \le \frac{1}{(p-2)C} \left(\frac{1}{z_0^{p-2}} - \frac{1}{z_1^{p-2}} \right) \le \frac{1}{(p-2)Cz_0^{p-2}}.$$
 (G.32)

Lemma 10 (Restated). Let p > 2. Given some C > 0, if for some z(t), the following holds

$$\frac{d}{dt}z \ge C(1-z), \forall t \in [0,T], \ z(0) = z_0, \ z(T) = z_1,$$
(G.33)

for some $0 < z_0 \le z_1 < 1$. Then the travel time T for z(t) to go from z_0 to z_1 satisfies:

$$T \le \frac{1}{C} \log \frac{1}{1 - z_1}$$
 (G.34)

Proof. We have

$$\int_{z_0}^{z_1} \frac{1}{C(1-z)} dz \ge \int_0^T dt \,, \tag{G.35}$$

thus

$$T \le \frac{1}{C} \left(\log \frac{1 - z_0}{1 - z_1} \right) \le \frac{1}{C} \log \frac{1}{1 - z_1}.$$
 (G.36)

Lemma 11 (Restated). Let p > 2. Condition on good event \mathcal{E}_{good} . Suppose the following is true at some point on the GF trajectory:

2634
2635 I.
$$c_{kj}(t) \ge 1 - 2C_a \log \frac{K}{\delta} \alpha^2, \ \forall k, j \in \mathcal{N}_k,$$

2636 2637 2.
$$\sum_{j \in \mathcal{N}_k} \| \boldsymbol{w}_j \|^2 \le 1 + C_w \log \frac{K}{\delta} \alpha^2, \ \forall k, j \in \mathcal{N}_k$$

2638 3.
$$\sum_{j \in \mathcal{N}_c} \| \boldsymbol{w}_j \|^2 = \tilde{o}(\alpha^2).$$
 2639

Then the following holds for every $1 \leq k \leq K$, $i \in \mathcal{I}_k$,

$$f^{(p)}(\boldsymbol{x}_i; \boldsymbol{\theta}) \leq \sum_{j \in \mathcal{N}_k} \|\boldsymbol{w}_j\|^2 \left(1 + 2^{p+2}C\sqrt{\log\frac{K^2N}{\delta}}\alpha^2\right) + 2KC\alpha^p$$

 $f^{(p)}(\boldsymbol{x}_i; \boldsymbol{\theta}) \geq \sum_{j \in \mathcal{N}_b} \|\boldsymbol{w}_j\|^2 \left(1 - 4pC\sqrt{\log \frac{K^2 N}{\delta}} \alpha^2\right) - 2KC\alpha^p.$

Proof. Our proof ignores terms related to neurons in \mathcal{N}_c as they only introduce a $\tilde{o}(\alpha^2)$ perturbation.

$$= \sum_{j=1}^{h} v_j \frac{\sigma^p(\langle \boldsymbol{w}_j, \boldsymbol{x}_i \rangle)}{\|\boldsymbol{w}_j\|^{p-1}}$$

$$= \sum_{j=1}^{h} \|\boldsymbol{w}_j\|^2 \sigma^p\left(\left\langle \frac{\boldsymbol{w}_j}{\|\boldsymbol{w}_j\|}, \boldsymbol{x}_i \right\rangle\right)$$

$$= \sum_{j \in \mathcal{N}_k} \|\boldsymbol{w}_j\|^2 \left(\left\langle \frac{\boldsymbol{w}_j}{\|\boldsymbol{w}_j\|}, \boldsymbol{x}_i \right\rangle\right)^p + \sum_{l \neq k} \sum_{j \in \mathcal{N}_l} \|\boldsymbol{w}_j\|^2 \sigma^p\left(\left\langle \frac{\boldsymbol{w}_j}{\|\boldsymbol{w}_j\|}, \boldsymbol{x}_i \right\rangle\right)$$

$$= \sum_{j \in \mathcal{N}_k} \|\boldsymbol{w}_j\|^2 \left(c_{kj} + \left\langle \frac{\boldsymbol{w}_j}{\|\boldsymbol{w}_j\|}, \boldsymbol{\varepsilon}_i \right\rangle\right)^p + \sum_{l \neq k} \sum_{j \in \mathcal{N}_l} \|\boldsymbol{w}_j\|^2 \sigma^p\left(c_{lj} + \left\langle \frac{\boldsymbol{w}_j}{\|\boldsymbol{w}_j\|}, \boldsymbol{\varepsilon}_i \right\rangle\right) \quad (G.37)$$

Upper bound:

 $f^{(p)}(\boldsymbol{x}_i;\boldsymbol{\theta})$

$$f^{(p)}(\boldsymbol{x}_{i};\boldsymbol{\theta}) = (G.37)$$

$$\leq \underbrace{\sum_{j \in \mathcal{N}_{k}} \|\boldsymbol{w}_{j}\|^{2} \left(c_{kj} + \left|\left\langle \frac{\boldsymbol{w}_{j}}{\|\boldsymbol{w}_{j}\|}, \boldsymbol{\varepsilon}_{i}\right\rangle\right|\right)^{p}}_{(a)} + \underbrace{\left|\sum_{l \neq k} \sum_{j \in \mathcal{N}_{l}} \|\boldsymbol{w}_{j}\|^{2} \sigma^{p} \left(c_{lj} + \left\langle \frac{\boldsymbol{w}_{j}}{\|\boldsymbol{w}_{j}\|}, \boldsymbol{\varepsilon}_{i}\right\rangle\right)\right|}_{(b)}$$

For the first term, we have

$$(a) \leq \sum_{j \in \mathcal{N}_k} \|\boldsymbol{w}_j\|^2 \left(1 + \|\boldsymbol{\varepsilon}_i\| \sqrt{1 - c_{kj}^2} + |\langle \boldsymbol{\mu}_k, \boldsymbol{\varepsilon}_i \rangle|\right)^p$$

$$\leq \sum_{j \in \mathcal{N}_k} \|\boldsymbol{w}_j\|^2 \left(1 + \|\boldsymbol{\varepsilon}_i\| \sqrt{2(1 - c_{kj})} + |\langle \boldsymbol{\mu}_k, \boldsymbol{\varepsilon}_i \rangle| \right)^p$$

$$\leq \sum_{j \in \mathcal{N}_k} \|\boldsymbol{w}_j\|^2 \left(1 + 2C\sqrt{\log\frac{K^2N}{\delta}}\alpha^2 + C\sqrt{\log\frac{K^2N}{\delta}}\frac{\alpha}{\sqrt{D}} \right)^p \\ \leq \sum_{j \in \mathcal{N}_k} \|\boldsymbol{w}_j\|^2 \left(1 + 2^{p+2}C\sqrt{\log\frac{K^2N}{\delta}}\alpha^2 \right) \,,$$

for sufficiently small α . For the second term, we have

$$(b) \leq 2 \sum_{l \neq k} \left(|c_{lj}| + \left| \left\langle \boldsymbol{\varepsilon}_{i}, \frac{\boldsymbol{w}_{j}}{\|\boldsymbol{w}_{j}\|} \right\rangle \right| \right)^{p} \\ \leq 2K \left(\sqrt{1 - c_{kj}^{2}} + \|\boldsymbol{\varepsilon}_{i}\| \sqrt{1 - c_{kj}^{2}} + |\langle \boldsymbol{\mu}_{k}, \boldsymbol{\varepsilon}_{i} \rangle| \right)^{p} \\ \leq 2K \left(\sqrt{2(1 - c_{kj})} + \|\boldsymbol{\varepsilon}_{i}\| \sqrt{(1 - c_{kj})} + |\langle \boldsymbol{\mu}_{k}, \boldsymbol{\varepsilon}_{i} \rangle| \right)^{p} \\ \leq 2K \left(C\alpha + C \sqrt{\log \frac{K^{2}N}{\delta}} \alpha^{2} + C \sqrt{\log \frac{K^{2}N}{\delta}} \frac{\alpha}{\sqrt{D}} \right)^{p} \leq 2KC\alpha^{p}.$$
(G.38)

Therefore

$$f^{(p)}(\boldsymbol{x}_i;\boldsymbol{\theta}) \le \sum_{j \in \mathcal{N}_k} \|\boldsymbol{w}_j\|^2 \left(1 + 2^{p+2}C\sqrt{\log\frac{K^2N}{\delta}}\alpha^2\right) + 2KC\alpha^p.$$
(G.39)

Lower bound:

$$f^{(p)}(\boldsymbol{x}_i; \boldsymbol{\theta})$$

$$\begin{array}{ll} 2700 \\ 2701 \\ 2702 \\ 2703 \\ 2& \sum_{j \in \mathcal{N}_{k}} \|\boldsymbol{w}_{j}\|^{2} \left(c_{kj} + \left|\left\langle \frac{\boldsymbol{w}_{j}}{\|\boldsymbol{w}_{j}\|}, \boldsymbol{\varepsilon}_{i}\right\rangle\right|\right)^{p} - \left|\sum_{l \neq k} \sum_{j \in \mathcal{N}_{l}} \|\boldsymbol{w}_{j}\|^{2} \sigma^{p} \left(c_{lj} + \left\langle \frac{\boldsymbol{w}_{j}}{\|\boldsymbol{w}_{j}\|}, \boldsymbol{\varepsilon}_{i}\right\rangle\right)\right)\right|. \\ 2704 \\ 2705 \\ 2706 \\ 2706 \\ 2706 \\ 2707 \\ 2708 \\ 26.38 \end{array}$$
For the first term, we have
$$\begin{array}{l} (a) \geq \sum_{j \in \mathcal{N}_{k}} \|\boldsymbol{w}_{j}\|^{2} \left(1 - \|\boldsymbol{\varepsilon}_{i}\|\sqrt{1 - c_{kj}^{2}} - |\langle\boldsymbol{\mu}_{k}, \boldsymbol{\varepsilon}_{i}\rangle|\right)^{p} \\ \geq \sum_{j \in \mathcal{N}_{k}} \|\boldsymbol{w}_{j}\|^{2} \left(1 - \|\boldsymbol{\varepsilon}_{i}\|\sqrt{2(1 - c_{kj})} - |\langle\boldsymbol{\mu}_{k}, \boldsymbol{\varepsilon}_{i}\rangle|\right)^{p} \\ \geq \sum_{j \in \mathcal{N}_{k}} \|\boldsymbol{w}_{j}\|^{2} \left(1 - 2C\sqrt{\log \frac{K^{2}N}{\delta}}\alpha^{2} - C\sqrt{\log \frac{K^{2}N}{\delta}}\frac{\alpha}{\sqrt{D}}\right)^{p} \\ \geq \sum_{j \in \mathcal{N}_{k}} \|\boldsymbol{w}_{j}\|^{2} \left(1 - 4pC\sqrt{\log \frac{K^{2}N}{\delta}}\alpha^{2}\right), \\ \text{for sufficiently small } \alpha. \text{ Therefore} \\ 2707 \\ 2718 \\ 2719 \\ 2721 \\ 2722 \\ 2722 \\ 2722 \\ 2724 \\ 2724 \\ 2724 \\ 2724 \\ 2724 \\ 2724 \\ 2725 \\ 2724 \\ 2724 \\ 2725 \\ 2724 \\ 2725 \\ 2724 \\ 2725 \\ 2726 \\ 2721 \\ 2726 \\ 2721 \\ 2726 \\ 2721 \\ 2726 \\ 2721 \\ 2726 \\ 2721 \\ 2726 \\ 2721 \\ 2726 \\ 2721 \\ 2726 \\ 2722 \\ 2724 \\ 2725 \\ 2726 \\ 2721 \\ 2726 \\ 2721 \\ 2726 \\ 2721 \\ 2726 \\ 2721 \\ 2726 \\ 2721 \\ 2726 \\ 2721 \\ 2726 \\ 2721 \\ 2726 \\ 2721 \\ 2726 \\ 2721 \\ 2726 \\ 2721 \\ 2726 \\ 2721 \\ 2726 \\ 2726 \\ 2721 \\ 2726 \\ 2726 \\ 2726 \\ 2726 \\ 2726 \\ 2726 \\ 2726 \\ 2726 \\ 2726 \\ 2726 \\ 2726 \\ 2726 \\ 2726 \\ 2726 \\ 2727 \\ 2726 \\ 2726 \\ 2727 \\ 2726 \\ 2726 \\ 2727 \\ 2726 \\ 2727 \\ 2726 \\ 2726 \\ 2727 \\ 2726 \\ 27$$

Lemma 12 (Restated). Let p > 2. Condition on good event \mathcal{E}_{good} . Suppose the following is true at some point on the GF trajectory:

2727 *I.*
$$c_{kj}(t) \ge 1 - 2C_a \log \frac{K}{\delta} \alpha^2, \ \forall k, j \in \mathcal{N}_k;$$

2729 2.
$$\sum_{j \in \mathcal{N}_k} \| \boldsymbol{w}_j \|^2 \le 1 + C_w \log \frac{K}{\delta} \alpha^2, \ \forall k,$$
2730

2731 3. $\sum_{j \in \mathcal{N}_c} \| \boldsymbol{w}_j \|^2 = \tilde{o}(\alpha^2).$

2734 2735

2737 2738

2739

Furthermore, suppose additionally that for some $k, j \in \mathcal{N}_k$:

$$1 - 2C_a \log \frac{K}{\delta} \alpha^2 \le c_{kj}(t) \le 1 - C_a \log \frac{K}{\delta} \alpha^2;$$

2736 Then the following holds for the same k, j,

$$\frac{d}{dt}c_{kj} \ge -CK\log\frac{K^2N}{\delta}\alpha^{\min\{p,4\}}.$$

2740 Proof. When $1 \le k \le K_1$, $j \in \mathcal{N}_k$ implies that $j \in \mathcal{N}_+$ thus $\operatorname{sign}(v_j) = 1$. We shall primarily focus 2741 on this case as the proof is nearly identical for $K_1 + 1 \le k \le K$. 2742 d

$$\begin{array}{ll}
2742 \\
2743 \\
2744 \\
2744 \\
2745 \\
2746 \\
2746 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\
2747 \\$$

$$\begin{array}{c} 2754\\ 2755\\ 2756\\ 2756\\ 2757\\ 2758 \end{array} + \underbrace{\frac{1}{N} \sum_{i \in \mathcal{I}_k} \left(\sum_{j \in \mathcal{N}_k} \|\boldsymbol{w}_j\|^2 - f^{(p)}(\boldsymbol{x}_i; \boldsymbol{\theta}) \right) p\left(\left\langle \boldsymbol{x}_i, \frac{\boldsymbol{w}_j}{\|\boldsymbol{w}_j\|} \right\rangle \right)^{p-1} \left(\left\langle \boldsymbol{\mu}_k, \boldsymbol{x}_i \right\rangle - \left\langle \boldsymbol{x}_i, \frac{\boldsymbol{w}_j}{\|\boldsymbol{w}_j\|} \right\rangle c_{kj} \right)}_{(b)} \right) }_{(b)}$$

$$+\underbrace{\frac{1}{N}\sum_{l\neq k}\sum_{i\in\mathcal{I}_{l}:\langle\boldsymbol{x}_{i},\boldsymbol{w}_{j}\rangle>0}(y_{i}-f^{(p)}(\boldsymbol{x}_{i};\boldsymbol{\theta}))p\left(\left\langle\left\langle\boldsymbol{x}_{i},\frac{\boldsymbol{w}_{j}}{\|\boldsymbol{w}_{j}\|}\right\rangle\right)^{p-1}\left(\left\langle\boldsymbol{\mu}_{k},\boldsymbol{x}_{i}\right\rangle-\left\langle\boldsymbol{x}_{i},\frac{\boldsymbol{w}_{j}}{\|\boldsymbol{w}_{j}\|}\right\rangle c_{kj}\right)}_{(c)}}_{(c)}$$
(G.41)

We deal with these terms one by one:

2766 Since $\sum_{j \in \mathcal{N}_k} \|\boldsymbol{w}_j\|^2 \le 1 + C\alpha^2$, for (a), there are two cases:

2768 1. When $1 - \sum_{j \in \mathcal{N}_k} \|\boldsymbol{w}_j\|^2 \ge 0$, Follow the same derivations from (G.12) to (G.14), we have

$$(a) = \left(1 - \sum_{j \in \mathcal{N}_{k}} \|\boldsymbol{w}_{j}\|^{2}\right) \frac{1}{N} \sum_{i \in \mathcal{I}_{k}} p\left(\left\langle\boldsymbol{x}_{i}, \frac{\boldsymbol{w}_{j}}{\|\boldsymbol{w}_{j}\|}\right\rangle\right)^{p-1} \left(\left\langle\boldsymbol{\mu}_{k}, \boldsymbol{x}_{i}\right\rangle - \left\langle\boldsymbol{x}_{i}, \frac{\boldsymbol{w}_{j}}{\|\boldsymbol{w}_{j}\|}\right\rangle c_{kj}\right)$$

$$\geq \left(1 - \sum_{j \in \mathcal{N}_{k}} \|\boldsymbol{w}_{j}\|^{2}\right) \left(pc_{kj}^{p-1}(1 - c_{kj}^{2}) - Cp^{2}\sqrt{\log\frac{K}{\delta}}\frac{\alpha}{\sqrt{N}} - 2^{p-1}p^{3}C^{2}\log\frac{K}{\delta}\alpha^{2} - o(\alpha^{2})\right)$$

$$\geq \left(1 - \sum_{j \in \mathcal{N}_{k}} \|\boldsymbol{w}_{j}\|^{2}\right) \left(p(1 - C\alpha^{2})^{p-1} \left(C\alpha^{2} - \frac{C^{2}}{4}\alpha^{4}\right) - Cp^{2}\sqrt{\log\frac{K}{\delta}}\frac{\alpha}{\sqrt{N}} - 2^{p-1}p^{3}C^{2}\log\frac{K}{\delta}\alpha^{2} - o(\alpha^{2})\right)$$

$$\geq 0,$$

for some choice of C and sufficiently small α .

2. When $-C\alpha^2 \leq 1 - \sum_{j \in \mathcal{N}_k} \|\boldsymbol{w}_j\|^2 \leq 0$, we have

$$(a) = \left(1 - \sum_{j \in \mathcal{N}_{k}} \|\boldsymbol{w}_{j}\|^{2}\right) \frac{1}{N} \sum_{i \in \mathcal{I}_{k}} p\left(\left\langle\boldsymbol{x}_{i}, \frac{\boldsymbol{w}_{j}}{\|\boldsymbol{w}_{j}\|}\right\rangle\right)^{p-1} \left(\left\langle\boldsymbol{\mu}_{k}, \boldsymbol{x}_{i}\right\rangle - \left\langle\boldsymbol{x}_{i}, \frac{\boldsymbol{w}_{j}}{\|\boldsymbol{w}_{j}\|}\right\rangle c_{kj}\right)$$

$$\geq -\left|1 - \sum_{j \in \mathcal{N}_{k}} \|\boldsymbol{w}_{j}\|^{2}\right| \frac{1}{N} \sum_{i \in \mathcal{I}_{k}} p\left(\left\langle\boldsymbol{x}_{i}, \frac{\boldsymbol{w}_{j}}{\|\boldsymbol{w}_{j}\|}\right\rangle\right)^{p-1} \left(1 - c_{kj}^{2} + |\langle\boldsymbol{\mu}_{k}, \boldsymbol{\varepsilon}_{i}\rangle| + \left|\left\langle\boldsymbol{\varepsilon}_{i}, \frac{\boldsymbol{w}_{j}}{\|\boldsymbol{w}_{j}\|}\right\rangle\right| c_{kj}\right)$$

$$\geq -\left|1 - \sum_{j \in \mathcal{N}_{k}} \|\boldsymbol{w}_{j}\|^{2}\right| p2^{p-1} \left(1 - c_{kj}^{2} + 2|\langle\boldsymbol{\mu}_{k}, \boldsymbol{\varepsilon}_{i}\rangle| + \|\boldsymbol{\varepsilon}_{i}\|\sqrt{1 - c_{kj}^{2}}c_{kj}\right)$$

$$\geq C\sqrt{\log \frac{K^{2}N}{\delta}} \alpha^{4}, \qquad (G.42)$$

Therefore, we always have

$$(a) \ge C\sqrt{\log \frac{K^2 N}{\delta}} \alpha^4 \,. \tag{G.43}$$

The second term (b) is easy: by Lemma 11, we know that $\left|\sum_{j\in\mathcal{N}_{k}}\|\boldsymbol{w}_{j}\|^{2} - f^{(p)}(\boldsymbol{x}_{i};\boldsymbol{\theta})\right| = \mathcal{O}(\sqrt{\log\frac{K^{2}N}{\delta}}\alpha^{2})$, then by the a similar derivation as in (G.42), we have (b) $= \frac{1}{N}\sum_{i\in\mathcal{I}_{k}}\left(\sum_{j\in\mathcal{N}_{k}}\|\boldsymbol{w}_{j}\|^{2} - f^{(p)}(\boldsymbol{x}_{i};\boldsymbol{\theta})\right) p\left(\left\langle\boldsymbol{x}_{i},\frac{\boldsymbol{w}_{j}}{\|\boldsymbol{w}_{j}\|}\right\rangle\right)^{p-1}\left(\left\langle\boldsymbol{\mu}_{k},\boldsymbol{x}_{i}\right\rangle - \left\langle\boldsymbol{x}_{i},\frac{\boldsymbol{w}_{j}}{\|\boldsymbol{w}_{j}\|}\right\rangle c_{kj}\right)$

$$\geq -\left|\sum_{j\in\mathcal{N}_{k}}\|\boldsymbol{w}_{j}\|^{2} - f^{(p)}(\boldsymbol{x}_{i};\boldsymbol{\theta})\right|\frac{1}{N}\sum_{i\in\mathcal{I}_{k}}p\left(\left\langle\boldsymbol{x}_{i},\frac{\boldsymbol{w}_{j}}{\|\boldsymbol{w}_{j}\|}\right\rangle\right)^{p-1}\left(1 - c_{kj}^{2} + |\langle\boldsymbol{\mu}_{k},\boldsymbol{\varepsilon}_{i}\rangle| + \left|\left\langle\boldsymbol{\varepsilon}_{i},\frac{\boldsymbol{w}_{j}}{\|\boldsymbol{w}_{j}\|}\right\rangle\right|c_{kj}\right)$$
$$\geq C\log\frac{K^{2}N}{\delta}\alpha^{4}, \tag{G.44}$$

For the last term, we have

$$(c) = \frac{1}{N} \sum_{l \neq k} \sum_{i \in \mathcal{I}_l: \langle \boldsymbol{x}_i, \boldsymbol{w}_j \rangle > 0} (y_i - f^{(p)}(\boldsymbol{x}_i; \boldsymbol{\theta})) p\left(\left\langle \boldsymbol{x}_i, \frac{\boldsymbol{w}_j}{\|\boldsymbol{w}_j\|} \right\rangle\right)^{p-1} \left(\left\langle \boldsymbol{\mu}_k, \boldsymbol{x}_i \right\rangle - \left\langle \boldsymbol{x}_i, \frac{\boldsymbol{w}_j}{\|\boldsymbol{w}_j\|} \right\rangle c_{kj}\right)$$

$$\geq -\frac{2}{N} \sum_{l \neq k} \sum_{i \in \mathcal{I}_l} p\left(c_{lj} + \left|\left\langle \boldsymbol{\varepsilon}_i, \frac{\boldsymbol{w}_j}{\|\boldsymbol{w}_j\|} \right\rangle\right|\right)^{p-1} \left(\left\langle \boldsymbol{\mu}_k, \boldsymbol{\varepsilon}_i \right\rangle + c_{lj}c_{kj} + \left|\left\langle \boldsymbol{\varepsilon}_i, \frac{\boldsymbol{w}_j}{\|\boldsymbol{w}_j\|} \right\rangle\right| c_{kj}\right)$$

$$\geq -\frac{2}{N} \sum_{l \neq k} \sum_{i \in \mathcal{I}_l} p\left(\sqrt{1 - c_{kj}^2} + \|\boldsymbol{\varepsilon}_i\| \sqrt{1 - c_{kj}^2}\right)^{p-1} \left(\left\langle \boldsymbol{\mu}_k, \boldsymbol{\varepsilon}_i \right\rangle + \sqrt{1 - c_{kj}^2}c_{kj} + \|\boldsymbol{\varepsilon}_i\| \sqrt{1 - c_{kj}^2}c_{kj}\right)$$

$$\geq -CK \sqrt{\log \frac{K^2 N}{\delta}} \alpha^p.$$

Finally, we can conclude that

$$\frac{d}{dt}c_{kj} \ge -CK\log\frac{K^2N}{\delta}\alpha^{\min\{p,4\}}.$$
(G.45)

Lemma 13 (Restated). Let p > 2. Condition on good event \mathcal{E}_{good} . Suppose the following is true at some point on the GF trajectory :

I.
$$c_{kj}(t) \ge 1 - 2C_a \log \frac{K}{\delta} \alpha^2, \ k, j \in \mathcal{N}_k;$$

2.
$$\sum_{j \in \mathcal{N}_k} \|\boldsymbol{w}_j\|^2 \le 1 + C_w \log \frac{K}{\delta} \alpha^2, \ \forall k;$$

2838 3.
$$\sum_{j \in \mathcal{N}_c} \| \boldsymbol{w}_j \|^2 = \tilde{o}(\alpha^2)$$

2840 Then the following holds for every $1 \le k \le K$,

$$\frac{d}{dt}\left(\sum_{j\in\mathcal{N}_k}\|\boldsymbol{w}_j\|^2\right) \leq 2\left(1-\sum_{j\in\mathcal{N}_k}\|\boldsymbol{w}_j\|^2+C\log\frac{K}{\delta}\alpha^2\right)\left(\sum_{j\in\mathcal{N}_k}\|\boldsymbol{w}_j\|^2\right)\,,$$

and

$$\frac{d}{dt}\left(\sum_{j\in\mathcal{N}_k}\|\boldsymbol{w}_j\|^2\right)\geq 2\left(1-\sum_{j\in\mathcal{N}_k}\|\boldsymbol{w}_j\|^2-C\log\frac{K}{\delta}\alpha^2\right)\left(\sum_{j\in\mathcal{N}_k}\|\boldsymbol{w}_j\|^2\right),$$

where C is some universal constant such that $C < C_w$.

Proof. When $1 \le k \le K_1$, $j \in \mathcal{N}_k$ implies that $j \in \mathcal{N}_+$ thus $sign(v_j) = 1$. We shall primarily focus on this case as the proof is nearly identical for $K_1 + 1 \le k \le K$. We start with (D.3):

$$\frac{d}{dt} \|\boldsymbol{w}_j\|^2 = \frac{2}{N} \left(\sum_{i: \langle \boldsymbol{x}_i, \boldsymbol{w}_j \rangle > 0} (y_i - f^{(p)}(\boldsymbol{x}_i; \boldsymbol{\theta})) \left(\left\langle \boldsymbol{x}_i, \frac{\boldsymbol{w}_j}{\|\boldsymbol{w}_j\|} \right\rangle \right)^p \right) \|\boldsymbol{w}_j\|^2$$

2858
2859
2860
2861
$$= 2\left(\underbrace{\frac{1}{N}\sum_{i\in\mathcal{I}_k:\langle \boldsymbol{x}_i,\boldsymbol{w}_j\rangle>0}(y_i - f^{(p)}(\boldsymbol{x}_i;\boldsymbol{\theta}))\left(\left\langle \boldsymbol{x}_i,\frac{\boldsymbol{w}_j}{\|\boldsymbol{w}_j\|}\right\rangle\right)^p}_{(a)}\right)$$

$$+\underbrace{\frac{1}{N}\sum_{l\neq k}\sum_{i\in\mathcal{I}_l:\langle \boldsymbol{x}_i,\boldsymbol{w}_j\rangle>0}(y_i-f^{(p)}(\boldsymbol{x}_i;\boldsymbol{\theta}))~\left(\left\langle \boldsymbol{x}_i,\frac{\boldsymbol{w}_j}{\|\boldsymbol{w}_j\|}\right\rangle\right)^p}_{:=\Gamma_1}}_{}_{}$$

For the first term, we have

$$\begin{aligned} & (a) = \frac{1}{N} \sum_{i \in \mathcal{I}_k: \langle \boldsymbol{x}_i, \boldsymbol{w}_j \rangle > 0} (y_i - f^{(p)}(\boldsymbol{x}_i; \boldsymbol{\theta})) \left(\left\langle \boldsymbol{x}_i, \frac{\boldsymbol{w}_j}{\|\boldsymbol{w}_j\|} \right\rangle \right)^p \\ & = \frac{1}{N} \sum_{i \in \mathcal{I}_k} (1 - f^{(p)}(\boldsymbol{x}_i; \boldsymbol{\theta})) \left(\left\langle \boldsymbol{x}_i, \frac{\boldsymbol{w}_j}{\|\boldsymbol{w}_j\|} \right\rangle \right)^p \\ & = \frac{1}{N} \sum_{i \in \mathcal{I}_k} \left(1 - \sum_{j \in \mathcal{N}_k} \|\boldsymbol{w}_j\|^2 \left(\left\langle \boldsymbol{x}_i, \frac{\boldsymbol{w}_j}{\|\boldsymbol{w}_j\|} \right\rangle \right)^p + \sum_{l \neq k} \sum_{j \in \mathcal{N}_l} \|\boldsymbol{w}_j\|^2 \sigma^p \left(\left\langle \boldsymbol{x}_i, \frac{\boldsymbol{w}_j}{\|\boldsymbol{w}_j\|} \right\rangle \right) \right) \left(\left\langle \boldsymbol{x}_i, \frac{\boldsymbol{w}_j}{\|\boldsymbol{w}_j\|} \right\rangle \right)^p \\ & = \frac{1}{N} \sum_{i \in \mathcal{I}_k} \left(1 - \sum_{j \in \mathcal{N}_k} \|\boldsymbol{w}_j\|^2 \left(\left\langle \boldsymbol{x}_i, \frac{\boldsymbol{w}_j}{\|\boldsymbol{w}_j\|} \right\rangle \right)^p \right) \left(\left\langle \boldsymbol{x}_i, \frac{\boldsymbol{w}_j}{\|\boldsymbol{w}_j\|} \right\rangle \right)^p + \underbrace{\frac{1}{N} \sum_{i \in \mathcal{I}_k} \sum_{l \neq k} \sum_{j \in \mathcal{N}_l} \|\boldsymbol{w}_j\|^2 \sigma^{2p} \left(\left\langle \boldsymbol{x}_i, \frac{\boldsymbol{w}_j}{\|\boldsymbol{w}_j\|} \right\rangle \right)}_{:=\Gamma_2} \\ & = \frac{1}{N} \sum_{i \in \mathcal{I}_k} \left(1 - \sum_{j \in \mathcal{N}_k} \|\boldsymbol{w}_j\|^2 \left(c_{kj} + \left\langle \boldsymbol{\varepsilon}_i, \frac{\boldsymbol{w}_j}{\|\boldsymbol{w}_j\|} \right\rangle \right)^p \right) \left(c_{kj} + \left\langle \boldsymbol{\varepsilon}_i, \frac{\boldsymbol{w}_j}{\|\boldsymbol{w}_j\|} \right\rangle \right)^p + \Gamma_2 . \end{aligned}$$

We shall focus on the first term. With the Taylor expansion

$$\left(c_{kj} + \left\langle \boldsymbol{\varepsilon}_{i}, \frac{\boldsymbol{w}_{j}}{\|\boldsymbol{w}_{j}\|} \right\rangle\right)^{p} = c_{kj}^{p} + pc_{kj}^{p-1} \left\langle \boldsymbol{\varepsilon}_{i}, \frac{\boldsymbol{w}_{j}}{\|\boldsymbol{w}_{j}\|} \right\rangle + R_{L} \left| \left\langle \boldsymbol{\varepsilon}_{i}, \frac{\boldsymbol{w}_{j}}{\|\boldsymbol{w}_{j}\|} \right\rangle \right|^{2}, \quad (G.46)$$

`

where $R_L = \frac{p(p-1)(c_{kj}+\zeta_L)^{p-2}}{2}$ and ζ_L between 0 and $\left|\left\langle \varepsilon_i, \frac{w_j}{\|w_j\|} \right\rangle\right|$ comes from the Lagrange residual. Clearly $|R_L| \leq 2^{p-2}p^2$. Then we have

$$\frac{1}{N}\sum_{i\in\mathcal{I}_{k}}\left(1-\sum_{j\in\mathcal{N}_{k}}\|\boldsymbol{w}_{j}\|^{2}\left(c_{kj}+\left\langle\boldsymbol{\varepsilon}_{i},\frac{\boldsymbol{w}_{j}}{\|\boldsymbol{w}_{j}\|}\right\rangle\right)^{p}\right)\left(c_{kj}+\left\langle\boldsymbol{\varepsilon}_{i},\frac{\boldsymbol{w}_{j}}{\|\boldsymbol{w}_{j}\|}\right\rangle\right)^{p}$$
$$=c_{kj}^{p}-\sum_{j\in\mathcal{N}_{k}}\|\boldsymbol{w}_{j}\|^{2}c_{kj}^{2p}$$
$$\left(pc_{kj}^{p-1}-2\sum_{j\in\mathcal{N}_{k}}\|\boldsymbol{w}_{j}\|^{2}c_{kj}^{2p-1}\right)\frac{1}{N}\sum_{i\in\mathcal{I}_{k}}\left\langle\boldsymbol{\varepsilon}_{i},\frac{\boldsymbol{w}_{j}}{\|\boldsymbol{w}_{j}\|}\right\rangle$$
$$\left(R_{L}-p^{2}c_{kj}^{2p-2}-2c_{kj}^{p}R_{L}\right)\frac{1}{N}\sum_{i\in\mathcal{I}_{k}}\left|\left\langle\boldsymbol{\varepsilon}_{i},\frac{\boldsymbol{w}_{j}}{\|\boldsymbol{w}_{j}\|}\right\rangle\right|^{2}+o\left(\left|\left\langle\boldsymbol{\varepsilon}_{i},\frac{\boldsymbol{w}_{j}}{\|\boldsymbol{w}_{j}\|}\right\rangle\right|^{2}\right)$$

Finally, we are ready to derive the upper and lower bound. For lower bound,

2916
2917
$$-o\left(\left|\left\langle \boldsymbol{\varepsilon}_{i}, \frac{\boldsymbol{w}_{j}}{\|\boldsymbol{w}_{j}\|} \right\rangle\right|^{2}\right) - |\Gamma_{1}| - |\Gamma_{2}|\right) \|\boldsymbol{w}_{j}\|^{2}$$
2918

$$2919 \\ \geq 2 \left(c_{kj}^{p} - \sum_{j \in \mathcal{N}_{k}} \| \boldsymbol{w}_{j} \|^{2} c_{kj}^{2p} - C \sqrt{\log \frac{K}{\delta}} \frac{\alpha}{\sqrt{N}} - C^{2} \log \frac{K}{\delta} \alpha^{2} - o(\alpha^{2}) - |\Gamma_{1}| - |\Gamma_{2}| \right) \| \boldsymbol{w}_{j} \|^{2}$$

$$2922 \\ 2923 \\ \geq 2 \left(\left(1 - \frac{C\alpha^{2}}{2} \right)^{p} - \sum_{j \in \mathcal{N}_{k}} \| \boldsymbol{w}_{j} \|^{2} - C \sqrt{\log \frac{K}{\delta}} \frac{\alpha}{\sqrt{N}} - C^{2} \log \frac{K}{\delta} \alpha^{2} - o(\alpha^{2}) - |\Gamma_{1}| - |\Gamma_{2}| \right) \| \boldsymbol{w}_{j} \|^{2}$$

$$2925 \\ 2926 \\ 2927 \\ \geq 2 \left(1 - p \frac{C\alpha^{2}}{2} - \sum_{j \in \mathcal{N}_{k}} \| \boldsymbol{w}_{j} \|^{2} - C \sqrt{\log \frac{K}{\delta}} \frac{\alpha}{\sqrt{N}} - C^{2} \log \frac{K}{\delta} \alpha^{2} - o(\alpha^{2}) - |\Gamma_{1}| - |\Gamma_{2}| \right) \| \boldsymbol{w}_{j} \|^{2}$$

$$2926 \\ \geq 2 \left(1 - p \frac{C\alpha^{2}}{2} - \sum_{j \in \mathcal{N}_{k}} \| \boldsymbol{w}_{j} \|^{2} - C \sqrt{\log \frac{K}{\delta}} \frac{\alpha}{\sqrt{N}} - C^{2} \log \frac{K}{\delta} \alpha^{2} - o(\alpha^{2}) - |\Gamma_{1}| - |\Gamma_{2}| \right) \| \boldsymbol{w}_{j} \|^{2}$$

It remains to bound these $|\Gamma_1|, |\Gamma_2|$. Indeed, we can find the following bound:

$$\begin{aligned} |\Gamma_{1}| &= \left| \frac{1}{N} \sum_{l \neq k} \sum_{i \in \mathcal{I}_{l}: \langle \boldsymbol{x}_{i}, \boldsymbol{w}_{j} \rangle > 0} (y_{i} - f^{(p)}(\boldsymbol{x}_{i}; \boldsymbol{\theta})) \left(\left\langle \boldsymbol{x}_{i}, \frac{\boldsymbol{w}_{j}}{\|\boldsymbol{w}_{j}\|} \right\rangle \right)^{p} \right| \\ &\leq \left| \frac{1}{N} \sum_{l \neq k} \sum_{i \in \mathcal{I}_{l}} |y_{i} - f^{(p)}(\boldsymbol{x}_{i}; \boldsymbol{\theta}))| \left(\left\langle \boldsymbol{x}_{i}, \frac{\boldsymbol{w}_{j}}{\|\boldsymbol{w}_{j}\|} \right\rangle \right)^{p} \right| \\ &\leq \left| \frac{2}{N} \sum_{l \neq k} \sum_{i \in \mathcal{I}_{l}} \sum_{l \neq k} \left(\left\langle \boldsymbol{x}_{i}, \frac{\boldsymbol{w}_{j}}{\|\boldsymbol{w}_{j}\|} \right\rangle \right)^{p} \right| \\ &\leq \left| \frac{2}{N} \sum_{i \in \mathcal{I}_{l}} \sum_{l \neq k} \left(\left\langle \boldsymbol{x}_{i}, \frac{\boldsymbol{w}_{j}}{\|\boldsymbol{w}_{j}\|} \right\rangle \right)^{p} \right| \\ &\leq \left| \frac{2}{N} \sum_{i \in \mathcal{I}_{l}} \sum_{l \neq k} \left(c_{lj} + \left\langle \boldsymbol{\varepsilon}_{i}, \frac{\boldsymbol{w}_{j}}{\|\boldsymbol{w}_{j}\|} \right\rangle \right)^{p} \right| \\ &\leq \left| \frac{2}{N} \sum_{i \in \mathcal{I}_{l}} \sum_{l \neq k} \left(\sqrt{1 - c_{kj}^{2}} + \|\boldsymbol{\varepsilon}_{i}\| \sqrt{1 - c_{kj}^{2}} + \langle \boldsymbol{\mu}_{k}, \boldsymbol{\varepsilon}_{i} \rangle \right)^{p} \right| \leq KC\alpha^{p}, \\ &|\Gamma_{2}| = \left| \frac{1}{N} \sum_{i \in \mathcal{I}_{k}} \sum_{l \neq k} \sum_{j \in \mathcal{N}_{l}} \|\boldsymbol{w}_{j}\|^{2} \sigma^{2p} \left(\left\langle \boldsymbol{x}_{i}, \frac{\boldsymbol{w}_{j}}{\|\boldsymbol{w}_{j}\|} \right\rangle \right) \right| \\ &\leq \left| \frac{2}{N} \sum_{i \in \mathcal{I}_{k}} \sum_{l \neq k} \left(\left\langle \boldsymbol{x}_{i}, \frac{\boldsymbol{w}_{j}}{\|\boldsymbol{w}_{j}\|} \right\rangle \right)^{2p} \right| \leq KC\alpha^{2p}. \end{aligned}$$

Therefore,

 $\frac{d}{dt} \| \boldsymbol{w}_j \|^2$

$$\frac{d}{dt} \|\boldsymbol{w}_j\|^2 \ge 2 \left(1 - \sum_{j \in \mathcal{N}_k} \|\boldsymbol{w}_j\|^2 - C \log \frac{K}{\delta} \alpha^2 \right) \|\boldsymbol{w}_j\|^2,$$

since when α is sufficiently small, the dominant term is of order α^2 .

Similarly, for the upper bound, we can have

$$= 2 \left((a) + \Gamma_1 \right) \| \boldsymbol{w}_j \|^2$$

$$\leq 2 \left(c_{kj}^p - \sum_{j \in \mathcal{N}_k} \| \boldsymbol{w}_j \|^2 c_{kj}^{2p} + C \sqrt{\log \frac{K}{\delta}} \frac{\alpha}{\sqrt{N}} + C^2 \log \frac{K}{\delta} \alpha^2 + o(\alpha^2) + |\Gamma_1| + |\Gamma_2| \right) \| \boldsymbol{w}_j \|^2$$

$$2970 \\
2971 \\
2972 \\
2973 \\
2973 \\
2974 \\
2974 \\
2975 \\
2976 \\
2976 \\
2976 \\
2976 \\
2976 \\
2976 \\
2976 \\
2976 \\
2976 \\
2976 \\
2976 \\
2976 \\
2976 \\
2976 \\
2976 \\
2976 \\
2976 \\
2976 \\
2976 \\
2976 \\
2976 \\
2976 \\
2976 \\
2976 \\
2976 \\
2976 \\
2976 \\
2976 \\
2976 \\
2976 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2970 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\
2070 \\$$

Lemma 14 (Restated). Consider the same assumptions as in Proposition 2. Given the t_1 in Proposition 2, the following holds $\forall 1 \le k \le K$:

$$\sum_{j\in\mathcal{N}_k} \|\boldsymbol{w}_j(t_1)\|^2 \ge \exp\left(-\frac{2p^{p+2}K}{p(p-2)\tilde{\Delta}_2\tilde{\Delta}_1^{p-2}}\right) W_{\min}^2 \epsilon^2.$$
(G.47)

Proof. The proof will be in two parts: first, we define, for each k,

$$t_{\text{aux}}^{(k)} := \inf \left\{ t : \min_{j \in \mathcal{N}_k} c_{kj}(t) \ge \frac{2}{3} \right\} \stackrel{\text{(By its definition)}}{\le} t_1 , \qquad (G.48)$$

and show that

$$\sum_{j \in \mathcal{N}_k} \|\boldsymbol{w}_j(t_{\text{aux}}^{(k)})\|^2 \ge \exp\left(-\frac{2p^{p+2}K}{p(p-2)\tilde{\Delta}_2\tilde{\Delta}_1^{p-2}}\right) \sum_{j \in \mathcal{N}_k} \|\boldsymbol{w}_j(0)\|^2.$$
(G.49)

Then we show that $\sum_{j \in \mathcal{N}_k} \| \boldsymbol{w}_j(t_1) \|^2$ is non-decreasing during $[t_{aux}^{(k)}, t_1]$.

Lower bound at $t_{aux}^{(k)}$: We shall focus on the case $1 \le k \le K_1$. In the proofs of Proposition 2, we have shown in (E.18) that when $t \le t_{aux}^{(k)} \le \overline{t}_1$, the following is true: $\forall j \in \mathcal{N}_k$

$$\frac{d}{dt}c_{kj} \ge \tilde{\Delta}_2 p c_{kj}^{p-1} \,, \tag{G.50}$$

3004 By Lemma 9, we have

$$t_{\text{aux}}^{(k)} = \inf\left\{t : c_{kj} \ge \frac{2}{3}\right\} \le \frac{1}{p(p-2)\tilde{\Delta}_2\tilde{\Delta}_1^{p-2}}.$$
 (G.51)

Now we are ready to lower bound $\sum_{j \in \mathcal{N}_k} \|\boldsymbol{w}_j(t_{aux}^{(k)})\|^2$: In the same way we derived (G.27), we can also obtain: for $t \le t_1$,

$$\frac{d}{dt} \|\boldsymbol{w}_j\|^2 \ge -2^{p+1} K (1 + 4\epsilon \sqrt{h} W_{\max}^2) \|\boldsymbol{w}_j\|^2 \ge -2^{p+2} K \|\boldsymbol{w}_j\|^2, \qquad (G.52)$$

thus

$$\frac{d}{dt}\sum_{j\in\mathcal{N}_k}\|\boldsymbol{w}_j\|^2 \ge -2^{p+2}K\sum_{j\in\mathcal{N}_k}\|\boldsymbol{w}_j\|^2.$$
(G.53)

3016 Finally, by Grönwall's inequality, we have

$$\sum_{j \in \mathcal{N}_{k}} \|\boldsymbol{w}_{j}(t_{\text{aux}}^{(k)})\|^{2} \ge \exp\left(-2p^{p+2}Kt_{\text{aux}}^{(k)}\right) \sum_{j \in \mathcal{N}_{k}} \|\boldsymbol{w}_{j}(0)\|^{2}$$

$$\sum_{j \in \mathcal{N}_{k}} \|\boldsymbol{w}_{j}(t_{\text{aux}}^{(k)})\|^{2} \ge \exp\left(-\frac{2p^{p+2}K}{p(p-2)\tilde{\Delta}_{2}\tilde{\Delta}_{1}^{p-2}}\right) W_{\min}^{2}\epsilon^{2}.$$

Norm is non-decreasing afterward The techniques we will be using here is similar to those used in proving previous lemma, so we describe the argument briefly.

3024 Suppose $1 \le k \le K$, we have the norm dynamics

$$\begin{split} \frac{d}{dt} \|\boldsymbol{w}_{j}\|^{2} \\ &= -\frac{2}{N} \left(\sum_{i:\langle \boldsymbol{x}_{i}, \boldsymbol{w}_{j} \rangle > 0} \nabla_{j} \ell_{i} \left(\left\langle \boldsymbol{x}_{i}, \frac{\boldsymbol{w}_{j}}{\|\boldsymbol{w}_{j}\|} \right\rangle \right)^{p} \right) \|\boldsymbol{w}_{j}\|^{2} \\ &= \frac{2}{N} \left(\sum_{i:\langle \boldsymbol{x}_{i}, \boldsymbol{w}_{j} \rangle > 0} y_{i} \left(\left\langle \boldsymbol{x}_{i}, \frac{\boldsymbol{w}_{j}}{\|\boldsymbol{w}_{j}\|} \right\rangle \right)^{p} \right) \|\boldsymbol{w}_{j}\|^{2} + \underbrace{\mathcal{O}(\epsilon)}_{\text{Recall how we handle } \Gamma_{1} \text{ in the proof of Lemma 6}} \\ &\geq \frac{2}{N} \left(\sum_{i\in\mathcal{I}_{k}} \left| \left\langle \boldsymbol{x}_{i}, \frac{\boldsymbol{w}_{j}}{\|\boldsymbol{w}_{j}\|} \right\rangle \right|^{p} - \sum_{l\neq k} \sum_{i\in\mathcal{I}_{l}} \left| \left\langle \boldsymbol{x}_{i}, \frac{\boldsymbol{w}_{j}}{\|\boldsymbol{w}_{j}\|} \right\rangle \right|^{p} \right) \|\boldsymbol{w}_{j}\|^{2} + \mathcal{O}(\epsilon) \\ &\geq \frac{2}{N} \left(\sum_{i\in\mathcal{I}_{k}} \underbrace{\left(c_{kj} + \left\langle \boldsymbol{x}_{i}, \frac{\boldsymbol{w}_{j}}{\|\boldsymbol{w}_{j}\|} \right\rangle \right)^{p}}_{\text{Taylor expansion, refer to (G.46)} - \sum_{l\neq k} \sum_{i\in\mathcal{I}_{l}} \underbrace{\left(|c_{lj}| + \left| \left\langle \boldsymbol{x}_{i}, \frac{\boldsymbol{w}_{j}}{\|\boldsymbol{w}_{j}\|} \right\rangle \right)^{p}}_{\text{Taylor expansion, refer to (G.46)} - \mathcal{O}\left(\alpha^{2} \right) \right) \|\boldsymbol{w}_{j}\|^{2} + \mathcal{O}(\epsilon) . \end{split}$$

When $c_{kj} \geq \frac{2}{3}$, we have

$$c_{kj}^{p} - \sum_{l \neq k} |c_{lj}|^{p} \ge c_{kj}^{p} - (1 - c_{kj}^{2})^{\frac{p}{2}} > 0, \qquad (G.54)$$

then for sufficiently small α and ϵ , we have $\frac{d}{dt} \| \boldsymbol{w}_j \|^2 \ge 0$. Then during $t_{aux}^{(k)} \le t \le t_1$, we have

$$\frac{d}{dt} \sum_{j \in \mathcal{N}_k} \|\boldsymbol{w}_j\|^2 \ge 0.$$
(G.55)

3056 The proof is finished.

Lemma 18. 15[Restated] Given some $0 < \Delta < \frac{1}{4}$, if for some z(t), the following holds

$$\frac{d}{dt}z \ge (1 - z - \Delta)z, \ z(0) = z_0, \ z(T) = z_1,$$
(G.56)

for some $0 < z_0 \le \frac{1}{4}$, and $z_0 \le z_1 < 1 - \Delta$. Then the travel time T for z(t) to go from z_0 to z_1 satisfies:

$$T \le 2\left(\log\frac{1}{1-z_1-\Delta} + \log\frac{1}{z_0}\right)$$
 (G.57)

Proof. We have

$$\int_{z_0}^{z_1} \frac{1}{(1-z-\Delta)z} dz \ge \int_0^T dt \,, \tag{G.58}$$

3070 thus

$$T \le \frac{1}{1 - \Delta} \left(\log \frac{1 - z_0 - \Delta}{1 - z_1 - \Delta} + \log \frac{z_1}{z_0} \right) \le 2 \left(\log \frac{1}{1 - z_1 - \Delta} + \log \frac{1}{z_0} \right) . \tag{G.59}$$

Lemma 16 (Restated). Condition on good event \mathcal{E}_{good} , we have

$$\sum_{j \in \mathcal{N}_c} \|\boldsymbol{w}_j(t)\|^2 = \tilde{o}(\alpha^2), \ \forall t \le T^* \,. \tag{G.60}$$

Proof. We deal with neurons with $sign(v_j) = +1$, the other case has a similar proof.

3080 If $j \in \mathcal{N}_c$, it means w_{j0} is initialized into the void region with $c_{kj}(0) < 0$ and $|c_{kj}(t)| = \Theta(1)$, for 3081 $1 \le k \le K_1$. Therefore, the inner product between $w_j(0)$ and a data point x_i from the k-th cluster is 3082 always negative, and this holds continuously as long as $c_{kj}(t) < 0$ and $|c_{kj}(t)| = \Theta(1)$.

We will show that

- 1. Until $t \leq t^*$, we still have $c_{kj}(t) < 0$ and $|c_{kj}(t)| = \Theta(1)$, thus none of the data in positive clusters activates w_j .
- 2. Then $c_{kj}(t) < 0$ and $|c_{kj}(t)| = \Theta(1)$ suggests that, $\sum_{j \in \mathcal{N}_c} \|\boldsymbol{w}_j\|^2$ has an at most $\mathcal{O}(\alpha^2)$ growth rate. And during $[t^*, T^*]$, with a slightly different argument, $\sum_{j \in \mathcal{N}_c} \|\boldsymbol{w}_j\|^2$ still has an at most $\mathcal{O}(\alpha^2)$ growth rate, thus continually stays at $\tilde{o}(\alpha^2)$.

The a more formal proof requires proof by contradiction, with previous lemmas we have proved, but the provided argument should easily be translated into a proof by contradiction.

First step: Given a $j \in \mathcal{N}_c \cup \mathcal{N}_+$ and $1 \le k \le K_1$, we have during $t \le t^*$,

$$\frac{d}{dt}c_{kj} = -\frac{1}{N}\sum_{i:\langle \boldsymbol{x}_i, \boldsymbol{w}_j \rangle > 0} \nabla_{\hat{y}} \ell_i \ p\left(\left\langle \boldsymbol{x}_i, \frac{\boldsymbol{w}_j}{\|\boldsymbol{w}_j\|} \right\rangle\right)^{p-1} \left(\langle \boldsymbol{\mu}_k, \boldsymbol{x}_i \rangle - \left\langle \boldsymbol{x}_i, \frac{\boldsymbol{w}_j}{\|\boldsymbol{w}_j\|} \right\rangle c_{kj}\right)$$

3101 3102

3103 3104

3108 3109 3110

3114311531163117

3120

3097

3084

3086

3087

3088

3089

3090 3091

3092

3093 3094

 $= -\frac{1}{N} \sum_{K_1+1 \le l \le K} \sum_{i \in \mathcal{I}_l: \langle \boldsymbol{x}_i, \boldsymbol{w}_j \rangle > 0} \nabla_{\hat{y}} \ell_i \, p\left(\left\langle \boldsymbol{x}_i, \frac{\boldsymbol{w}_j}{\|\boldsymbol{w}_j\|} \right\rangle\right)^{p-1} \left(\underbrace{\langle \boldsymbol{\mu}_k, \boldsymbol{x}_i \rangle}_{=\mathcal{O}(\frac{\alpha}{\sqrt{D}})} - \left\langle \boldsymbol{x}_i, \frac{\boldsymbol{w}_j}{\|\boldsymbol{w}_j\|} \right\rangle c_{kj}\right)$ $= -\frac{1}{N} \sum_{K_1+1 \le l \le K} \sum_{i \in \mathcal{I}_l: \langle \boldsymbol{x}_i, \boldsymbol{w}_j \rangle > 0} \nabla_{\hat{y}} \ell_i \, p\left(\left\langle \boldsymbol{x}_i, \frac{\boldsymbol{w}_j}{\|\boldsymbol{w}_j\|} \right\rangle\right)^p \underbrace{c_{kj}}_{<0} + \mathcal{O}\left(\frac{\alpha}{\sqrt{D}}\right).$

3105 3106 3106 3107 Since $\nabla_{\hat{y}}\ell_i$ is either < 0 (during alignment phase) or = $\mathcal{O}(\alpha^2)$ (after norm growth). Then we always have $\frac{d}{dt}c_{kj} = \mathcal{O}(\alpha^2)$. Therefore, $\forall t \leq t^*$

$$c_{kj}(t) \le c_{kj}(0) + t \cdot \mathcal{O}(\alpha^2) \le c_{kj}(0) + t^* \mathcal{O}(\alpha^2) = c_{kj}(0) + \mathcal{O}\left(\alpha^2 \log \frac{1}{\alpha}\right), \qquad (G.61)$$

thus, we still have $c_{kj}(t) < 0$ and $|c_{kj}(t)| = \Theta(1)$.

Second step: During $[0, t^*]$, since none of the data in positive clusters activates w_i , we have

$$\frac{d}{dt} \|\boldsymbol{w}_j\|^2 = -2 \frac{1}{N} \left(\sum_{i: \langle \boldsymbol{x}_i, \boldsymbol{w}_j \rangle > 0} \nabla_{\hat{y}} \ell_i \left(\left\langle \boldsymbol{x}_i, \frac{\boldsymbol{w}_j}{\|\boldsymbol{w}_j\|} \right\rangle \right)^p \right) \|\boldsymbol{w}_j\|^2$$

$$= -2 \frac{1}{N} \left(\sum_{K_1 + 1 \le l \le K} \sum_{i \in \mathcal{I}_l : \langle \boldsymbol{x}_i, \boldsymbol{w}_j \rangle > 0} \nabla_{\hat{y}} \ell_i \left(\left\langle \boldsymbol{x}_i, \frac{\boldsymbol{w}_j}{\|\boldsymbol{w}_j\|} \right\rangle \right)^p \right) \|\boldsymbol{w}_j\|^2.$$

3121 3122 Since $\nabla_{\hat{y}}\ell_i$ is either < 0 (during alignment phase) or = $\mathcal{O}(\alpha^2)$ (after norm growth). We have 3123 $\frac{d}{dt} \|\boldsymbol{w}_j\|^2 = \mathcal{O}(\alpha^2) \cdot \|\boldsymbol{w}_j\|^2$.

3124 During $[t^*, T^*]$, we have $\nabla_{\hat{y}} \ell_i = \mathcal{O}(\alpha^2)$ for all i (as the consequence of $\left|1 - \sum_{j \in \mathcal{N}_k} \|\boldsymbol{w}_j\|^2\right| = \mathcal{O}(\alpha^2)$ and Lemma 11). Therefore we still have $\frac{d}{dt} \|\boldsymbol{w}_j\|^2 = \mathcal{O}(\alpha^2) \cdot \|\boldsymbol{w}_j\|^2$.

Then we have $\forall t \leq T^*$, $\frac{d}{dt} \sum_{j \in \mathcal{N}_c} \|\boldsymbol{w}_j(t)\|^2 \leq \mathcal{O}\left(\exp(\alpha^2 T^*)\right) \sum_{j \in \mathcal{N}_c} \|\boldsymbol{w}_j(0)\|^2 \leq \mathcal{O}(1) \sum_{j \in \mathcal{N}_c} \|\boldsymbol{w}_j(0)\|^2 \leq \mathcal{O}(\epsilon^2) = \tilde{o}(\alpha^2).$ (G.62) **Lemma 17** (Restated). If the neurons $\{w_j\}_{j=1}^h$ satisfies the following for some $0 \le \delta \le 1$ and $\nu, \zeta > 0$: • $\max_k \max_{j \in \mathcal{N}_k} c_{kj}(t) \ge 1 - \delta;$ • $\left|1-\sum_{j\in\mathcal{N}_k}\|\boldsymbol{w}_j\|^2\right|\leq\nu;$ • $\sum_{j \in \mathcal{N}^c} \|\boldsymbol{w}_j\|^2 \leq \zeta$, then $\sup_{\boldsymbol{x} \in \mathbb{S}^{D-1}} \left| f^{(p)}(\boldsymbol{x}; \boldsymbol{\theta}) - F^{(p)}(\boldsymbol{x}) \right| \le K(1+\nu)(2^p-1)2\delta + K\nu + \zeta$ Proof. $f^{(p)}(\boldsymbol{x};\boldsymbol{\theta})$ (G.63) $= \sum_{i=1}^{h} v_j \frac{\sigma^p(\langle \boldsymbol{w}_j, \boldsymbol{x} \rangle)}{\|\boldsymbol{w}_j\|^{p-1}}$ $= \sum_{i=1}^{h} \operatorname{sign}(v_j) \|\boldsymbol{w}_j\|^2 \frac{\sigma^p(\langle \boldsymbol{w}_j, \boldsymbol{x} \rangle)}{\|\boldsymbol{w}_j\|^p}$ $= \sum_{i=1}^{n} \operatorname{sign}(v_{j}) \|\boldsymbol{w}_{j}\|^{2} \sigma^{p} \left(\left\langle \frac{\boldsymbol{w}_{j}}{\|\boldsymbol{w}_{j}\|}, \boldsymbol{x} \right\rangle \right)$ $\mathcal{L} = \sum_{1 \leq k \leq K_1} \sum_{j \in \mathcal{N}_k} \| oldsymbol{w}_j \|^2 \sigma^p \left(\left\langle \frac{oldsymbol{w}_j}{\|oldsymbol{w}_j\|}, oldsymbol{x}
ight
angle
ight) - \sum_{K_1 + 1 \leq k \leq K} \sum_{j \in \mathcal{N}_k} \| oldsymbol{w}_j \|^2 \sigma^p \left(\left\langle \frac{oldsymbol{w}_j}{\|oldsymbol{w}_j\|}, oldsymbol{x}
ight
angle
ight)$ $+ \sum_{j \in \mathcal{N}^{c}} \operatorname{sign}(v_{j}) \|\boldsymbol{w}_{j}\|^{2} \sigma^{p} \left(\left\langle \frac{\boldsymbol{w}_{j}}{\|\boldsymbol{w}_{j}\|}, \boldsymbol{x} \right\rangle \right)$ (G.64)

For the first term, we have $\forall x \in \mathbb{S}^{D-1}$

$$\begin{aligned} & |\sum_{1 \le k \le K_1} \sum_{j \in \mathcal{N}_k} \|w_j\|^2 \sigma^p \left(\left\langle \frac{w_j}{\|w_j\|}, x \right\rangle \right) - \sum_{1 \le k \le K_1} \sigma^p (\langle \mu_k, x \rangle) \right| \\ & \leq \sum_{1 \le k \le K_1} \left| \sum_{j \in \mathcal{N}_k} \|w_j\|^2 \sigma^p \left(\left\langle \frac{w_j}{\|w_j\|} - \mu_k + \mu_k, x \right\rangle \right) - \sigma^p (\langle \mu_k, x \rangle) \right| \\ & \leq \sum_{1 \le k \le K_1} \left| \sum_{j \in \mathcal{N}_k} \|w_j\|^2 \sigma^p \left(\langle \mu_k, x \rangle + \left\| \frac{w_j}{\|w_j\|} - \mu_k \right\| \right) - \sigma^p (\langle \mu_k, x \rangle) \right| \\ & \leq \sum_{1 \le k \le K_1} \left| \sum_{j \in \mathcal{N}_k} \|w_j\|^2 \sigma^p \left(\langle \mu_k, x \rangle + 2(1 - c_{kj}) \right) - \sigma^p (\langle \mu_k, x \rangle) \right| \\ & = \sum_{1 \le k \le K_1} \left| \sum_{j \in \mathcal{N}_k} \|w_j\|^2 \sigma^p \left(\langle \mu_k, x \rangle + 2(1 - c_{kj}) \right) - \sigma^p (\langle \mu_k, x \rangle) \right| \\ & \leq \sum_{1 \le k \le K_1} \left| \sum_{j \in \mathcal{N}_k} \|w_j\|^2 \sigma^p \left(\langle \mu_k, x \rangle + 2(1 - c_{kj}) \right) - \sigma^p (\langle \mu_k, x \rangle) \right| \\ & + \sum_{1 \le k \le K_1} \left| \sum_{j \in \mathcal{N}_k} \|w_j\|^2 \sigma^p (\langle \mu_k, x \rangle + 2(1 - c_{kj})) - \sigma^p (\langle \mu_k, x \rangle) \right| \\ & \leq \sum_{1 \le k \le K_1} (1 + \nu) |\sigma^p \left(\langle \mu_k, x \rangle + 2(1 - c_{kj}) \right) - \sigma^p (\langle \mu_k, x \rangle) | + \sum_{1 \le k \le K_1} \nu |\sigma^p (\langle \mu_k, x \rangle) + 2(1 - c_{kj})) - \sigma^p (\langle \mu_k, x \rangle) | + K_1 \nu \end{aligned}$$

 $\leq K_1(1+\nu)(2^p-1)2\delta + K_1\nu$,

where the last inequality is due to the following derivation (notice that ReLU $\sigma(z)$ is non-decreasing in z, and polynomial z^p is non-decreasing for z > 0)

3190
3190

$$|\sigma^{p}(\langle \boldsymbol{\mu}_{k}, \boldsymbol{x} \rangle + 2(1 - c_{kj})) - \sigma^{p}(\langle \boldsymbol{\mu}_{k}, \boldsymbol{x} \rangle)|$$

 $=\sigma^{p}(\langle \boldsymbol{\mu}_{k}, \boldsymbol{x} \rangle + 2(1 - c_{kj})) - (\langle \boldsymbol{\mu}_{k}, \boldsymbol{x} \rangle)$

$$\leq (1+2\delta)^p - 1 \leq (2^p - 1)2\delta.$$

3194 Similarly, for the second term, we have $\forall x \in \mathbb{S}^{D-1}$

$$\left| \sum_{K_1+1 \le k \le K} \sum_{j \in \mathcal{N}_k} \|\boldsymbol{w}_j\|^2 \sigma^p \left(\left\langle \frac{\boldsymbol{w}_j}{\|\boldsymbol{w}_j\|}, \boldsymbol{x} \right\rangle \right) - \sum_{K_1+1 \le k \le K} \sigma^p(\langle \boldsymbol{\mu}_k, \boldsymbol{x} \rangle) \right|$$

$$\leq K_2(1+\nu)(2^p-1)2\delta + K_2\nu$$

Lastly, for the third term, we have

$$\left|\sum_{j\in\mathcal{N}^c}\operatorname{sign}(v_j)\|\boldsymbol{w}_j\|^2\sigma^p\left(\left\langle\frac{\boldsymbol{w}_j}{\|\boldsymbol{w}_j\|},\boldsymbol{x}\right\rangle\right)\right|\leq \sum_{j\in\mathcal{N}^c}\|\boldsymbol{w}_j\|^2\leq \zeta$$

Therefore, for any $\boldsymbol{x} \in \mathbb{S}^{D-1}$, we have

$$\begin{aligned} \left| f^{(p)}(\boldsymbol{x};\boldsymbol{\theta}) - F^{(p)}(\boldsymbol{x}) \right| \\ \leq \left| \sum_{1 \leq k \leq K_1} \sum_{j \in \mathcal{N}_k} \|\boldsymbol{w}_j\|^2 \sigma^p \left(\left\langle \frac{\boldsymbol{w}_j}{\|\boldsymbol{w}_j\|}, \boldsymbol{x} \right\rangle \right) - \sum_{1 \leq k \leq K_1} \sigma^p(\langle \boldsymbol{\mu}_k, \boldsymbol{x} \rangle) \right| \\ + \left| \sum_{K_1 + 1 \leq k \leq K} \sum_{j \in \mathcal{N}_k} \|\boldsymbol{w}_j\|^2 \sigma^p \left(\left\langle \frac{\boldsymbol{w}_j}{\|\boldsymbol{w}_j\|}, \boldsymbol{x} \right\rangle \right) - \sum_{K_1 + 1 \leq k \leq K} \sigma^p(\langle \boldsymbol{\mu}_k, \boldsymbol{x} \rangle) \right| \\ + \left| \sum_{j \in \mathcal{N}^c} \operatorname{sign}(v_j) \|\boldsymbol{w}_j\|^2 \sigma^p \left(\left\langle \frac{\boldsymbol{w}_j}{\|\boldsymbol{w}_j\|}, \boldsymbol{x} \right\rangle \right) \right| \\ \leq K(1 + \nu)(2^p - 1) 2\delta + K\nu + \zeta \end{aligned}$$

-	-