

---

# Can Implicit Bias Imply Adversarial Robustness?

---

Hancheng Min<sup>1</sup> René Vidal<sup>1</sup>

## Abstract

The implicit bias of gradient-based training algorithms has been considered mostly beneficial as it leads to trained networks that often generalize well. However, Frei et al. (2023) show that such implicit bias can harm adversarial robustness. Specifically, when the data consists of clusters with small inter-cluster correlation, a shallow (two-layer) ReLU network trained by gradient flow generalizes well, but it is not robust to adversarial attacks of small radius, despite the existence of a much more robust classifier that can be explicitly constructed from a shallow network. In this paper, we extend recent analyses of neuron alignment to show that a shallow network with a polynomial ReLU activation (pReLU) trained by gradient flow not only generalizes well but is also robust to adversarial attacks. Our results highlight the importance of the interplay between data structure and architecture design in the implicit bias and robustness of trained networks.

## 1. Introduction

Behind the success of deep neural networks in many application domains lies their vulnerability to *adversarial attacks*, i.e., small and human-imperceptible perturbations to the input data. Such a phenomenon was observed in the seminal paper of Szegedy et al. (2014) and has motivated a large body of work on building defenses against such 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).

However, many defense strategies have been shown to fail against new adaptive attacks (Athalye et al., 2018; Carlini et al., 2019), and understanding these failures seems to be a fundamental challenge. For example, Fawzi et al.

<sup>1</sup>Department of Electrical and Systems Engineering, University of Pennsylvania. Correspondence to: Hancheng Min <hanchengmin@seas.upenn.edu>.

*Proceedings of the 41<sup>st</sup> International Conference on Machine Learning*, Vienna, Austria. PMLR 235, 2024. Copyright 2024 by the author(s).

(2018); Dohmatob (2019); Shafahi et al. (2018) show the non-existence of robust classifiers for certain data distributions. Recently, Pal et al. (2023) show that having a data distribution that concentrates within a small-volume subset of the ambient space is necessary for the existence of a robust classifier. These results highlight the importance of understanding and exploiting data structure in the process of finding classifiers with certified robustness, yet almost none of the existing defense strategies do so.

Besides data distribution, additional issues arise from the training algorithms. Vardi et al. (2022); Frei et al. (2023) have shown that when the data consists of clusters with small inter-cluster correlation, a shallow (two-layer) ReLU network trained by gradient flow generalizes well while failing to be robust against adversarial attacks of small radius, despite the existence of a much more robust classifier that can be explicitly constructed from a shallow network. This study unveils new challenges in the search for robust classifiers: Even if we know a robust classifier exists for certain data distribution, the *implicit bias* from our training algorithm (the choice of network architecture, optimization algorithm, etc.) may prevent us from finding it.

**Paper contributions.** In this paper, we show that under the same setting studied in Frei et al. (2023), the implicit bias of gradient flow that leads to non-robust networks can be altered to favor robust networks by modifying the ReLU activation. Specifically, we consider a data distribution consisting of a mixture of  $K$  Gaussians, referred to as *subclasses*, which have small inter-subclass correlation and are grouped into two *superclasses/classes*. When training a shallow network to predict the class membership of a given data, we show, partially with formal theorems, and partially with conjectures validated by empirical experiments, that:

- If the activation is a ReLU, neurons (rows of first-layer weight matrix) within the network tend to learn only the average direction of each class, leading to a classifier that generalizes well on clean data, but is vulnerable to an adversarial attack with  $l_2$  radius  $\mathcal{O}\left(\frac{1}{\sqrt{K}}\right)$ , i.e. the trained network is non-robust with many subclasses. This leads to a new neural alignment perspective on the nonrobustness of ReLU networks identified by Frei et al. (2023).
- If the activation is replaced by a novel polynomial ReLU

activation, proposed based on recent advances in understanding the neuron alignment in shallow networks, neurons tend to learn the direction of each subclass center, leading to a classifier that generalizes well on clean data and can sustain any adversarial attack with  $\mathcal{O}(1)$  radius.

Our analysis **(1)** highlights the importance of the interplay between data structure and network architecture in determining the robustness of the trained network, **(2)** explains how the implicit bias (regularization) of training a ReLU network fails to exploit the data structure and leads to non-robust networks, and **(3)** shows how the issue is resolved by using a polynomial ReLU activation function. Moreover, numerical experiments on real datasets show that shallow networks with our generalized ReLU activation functions are much more robust than those with ReLU activation.

**Relation to existing analysis on implicit bias of neural networks.** Our discussion is theory-centric. It builds upon the analysis of the implicit bias of training algorithms, but it is significantly different from prior theoretical analyses. The implicit bias of training algorithms has been studied extensively over past years for various architectures, including linear networks (Saxe et al., 2014; Gunasekar et al., 2017; Ji & Telgarsky, 2019; Woodworth et al., 2020; Min et al., 2021; Stöger & Soltanolkotabi, 2021; Jacot et al., 2021; Wang & Jacot, 2023), and nonlinear networks (Lyu & Li, 2019; Ji & Telgarsky, 2020; Maennel et al., 2018; Chizat & Bach, 2020; Boursier et al., 2022; Min et al., 2024; Wang & Ma, 2023; Frei et al., 2022; 2023; Kumar & Haupt, 2024; Abbe et al., 2023; Tarzanagh et al., 2023), and from different perspectives, including max-margin (Lyu & Li, 2019; Chizat & Bach, 2020; Tarzanagh et al., 2023), min-norm (Gunasekar et al., 2017; Min et al., 2021), sparsity/low-rankness (Saxe et al., 2014; Woodworth et al., 2020; Wang & Jacot, 2023; Abbe et al., 2023), and alignment (Maennel et al., 2018; Kumar & Haupt, 2024). While most works consider these implicit biases toward simple networks beneficial and prominent explanations for the success of neural networks in practice, and some existing work (Faghri et al., 2021) has even shown that such biases help robustness in the cases of linear regression; few works (Frei et al., 2023; Boursier & Flammarion, 2024) discuss the potential harm caused by such biases. Our work takes one step further by proposing fixes to the harms identified by prior work (Frei et al., 2023), which sheds light on the potential of using deep learning theory to not only understand but to also improve neural networks in practice.

**Paper organization.** Our main formal results regarding adversarial robustness will be shown (in Section 3) for explicitly constructed classifiers that can be realized by shallow networks with different activation functions. Then we connect these robustness results to those of actual trained

shallow networks by a conjecture that shallow networks, when initialized properly, can learn such constructions via gradient flow training. The rationale behind this conjecture is carefully explained in Section 4, together with a preliminary theoretical analysis that supports this conjecture. Lastly, we empirically verify our conjecture in Section 5, followed by numerical experiments showing that our proposed new activation improves the robustness of shallow networks on real datasets.

**Notation:** We denote the inner product between vectors  $\mathbf{x}$  and  $\mathbf{y}$  by  $\langle \mathbf{x}, \mathbf{y} \rangle = \mathbf{x}^\top \mathbf{y}$ , and the cosine of the angle between them as  $\cos(\mathbf{x}, \mathbf{y}) = \langle \frac{\mathbf{x}}{\|\mathbf{x}\|}, \frac{\mathbf{y}}{\|\mathbf{y}\|} \rangle$ . For an  $n \times m$  matrix  $\mathbf{A}$ , we let  $\|\mathbf{A}\|$  and  $\|\mathbf{A}\|_F$  denote the spectral norm and Frobenius norm of  $\mathbf{A}$ , 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  $\text{Unif}(S)$  denote the uniform distribution over a set  $S$ . Lastly, we let  $[N]$  denote the integer set  $\{1, \dots, N\}$ .

## 2. Problem Settings

We consider a binary classification problem where within each class, the data consists of subclasses with small inter-subclass correlations, formally defined as follows.

**Data distribution.** We assume  $(X, Y, Z) \sim \mathcal{D}$ , where a sample  $(\mathbf{x}, y, z) \in \mathbb{R}^D \times \{-1, 1\} \times \mathbb{N}_+$  drawn from  $\mathcal{D}$  consists of the *observed data* and *class label*  $(\mathbf{x}, y)$ , and a latent (unobserved) variable  $z$  denoting the *subclass membership*. We denote the marginal distribution of the observed part  $(X, Y)$  as  $\mathcal{D}_{X,Y}$ . In particular, given  $1 \leq K_1 < K$ , we let

- $Z \sim \mathcal{D}_Z = \text{Unif}([K]); Y = \mathbb{1}_{\{Z \leq K_1\}} - \mathbb{1}_{\{Z > K_1\}}$ ;
- $(X \mid Z = k) \sim \mathcal{N}\left(\boldsymbol{\mu}_k, \frac{\alpha^2}{D} \mathbf{I}\right)$ ,  $\forall 1 \leq k \leq K$ , where  $\alpha > 0$  is the *intra-subclass variance*, and  $\{\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_K\}$  are *subclass centers* that forms an orthonormal basis of a  $K$ -dimensional subspace in  $\mathbb{R}^D$ .

This data distribution has  $K$  subclasses with small inter-subclass correlation if  $\alpha$  is small (since the subclass centers are orthonormal to each other), and the observed labels only reveal the superclass/class membership: the first  $K_1$  subclasses belong to the positive class ( $y = +1$ ) and the remaining  $K_2 := K - K_1$  belong to the negative class ( $y = -1$ ). One such a dataset is depicted in Figure 1.

**Classification task and robustness of a classifier.** The learning task we consider is to find a binary classifier  $f(\mathbf{x}) : \mathbb{R}^D \rightarrow \mathbb{R}$  such that given a randomly drawn  $(\mathbf{x}, y)$ ,  $\text{sign}(f(\mathbf{x}))$  predicts the correct class label  $y$ , i.e.,  $f(\mathbf{x})y > 0$ , with high probability. Moreover, we are interested in the robustness of  $f(\mathbf{x})$  against additive adversarial attacks/perturbations on  $\mathbf{x}$  with  $l_2$ -norm bounded by

some constant  $r > 0$  (called the *attack radius*). Specifically, given a randomly drawn  $(\mathbf{x}, y)$  from  $\mathcal{D}_{X,Y}$ , we would like  $\inf_{\|\mathbf{d}\|=1} f(\mathbf{x} + r\mathbf{d})y > 0$  to hold with high probability (formal statement in the next section.) for  $r$  as large as possible.

**Gradient flow training.** A candidate classifier can be realized from a neural network  $f(\mathbf{x}; \boldsymbol{\theta})$  parametrized by its networks weights  $\boldsymbol{\theta}$ . Given a size- $n$  training dataset  $\{(\mathbf{x}_i, y_i), i = 1, \dots, n\}$ , where each training sample is randomly drawn from  $\mathcal{D}_{X,Y}$ , we define a loss function

$$\mathcal{L}(\boldsymbol{\theta}; \{\mathbf{x}_i, y_i\}_{i=1}^n) = \sum_{i=1}^n \ell(y_i, f(\mathbf{x}_i; \boldsymbol{\theta})), \quad (1)$$

where  $\ell(y, \hat{y}) : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$  is either the exponential loss  $\ell(y, \hat{y}) = \exp(-y\hat{y})$  or the logistic loss  $\ell(y, \hat{y}) = \log(1 + \exp(-y\hat{y}))$ . We train the network using *gradient flow* (GF),  $\dot{\boldsymbol{\theta}}(t) \in \partial \mathcal{L}(\boldsymbol{\theta}(t))$ , where  $\partial \mathcal{L}$  denotes the *Clarke subdifferential* (Clarke, 1990) w.r.t.  $\boldsymbol{\theta}$ . With proper initialization  $\boldsymbol{\theta}(0)$ , one expects the trained network weights  $\boldsymbol{\theta}(T)$  (for some large  $T > 0$ ) to be close to some minimizer of  $\mathcal{L}(\boldsymbol{\theta})$  and  $f(\mathbf{x}; \boldsymbol{\theta}(T))$  be a good classifier for  $\mathcal{D}_{X,Y}$ .

**Shallow networks with pReLU activation.** We specifically study the following shallow (two-layer) network, parametrized by  $\boldsymbol{\theta} := \{(\mathbf{w}_j, v_j) \in \mathbb{R}^D \times \mathbb{R}, j = 1, \dots, h\}$ , as a candidate classifier:

$$f_p(\mathbf{x}; \{\mathbf{w}_j, v_j\}_{j=1}^h) = \sum_{j=1}^h v_j \frac{[\sigma(\langle \mathbf{x}, \mathbf{w}_j \rangle)]^p}{\|\mathbf{w}_j\|^{p-1}}, \quad (\text{Shallow pReLU Network})$$

where  $\sigma(\cdot) = \text{ReLU}(\cdot) = \max\{\cdot, 0\}$  and  $p \geq 1$ . When  $p = 1$ , this is exactly a shallow ReLU network. When  $p \geq 1$ ,  $f_p$  can be loosely viewed as a shallow network with a polynomial ReLU activation and a special form of weight normalization (Salimans & Kingma, 2016) on each  $\mathbf{w}_j$ . One of the most important reasons for this particular generalization of the ReLU activation is the “extra penalty” on angle separation between the input  $\mathbf{x}$  and *neurons*  $\mathbf{w}_j$ s. To see this, assume  $\|\mathbf{x}\| = 1$ , then

$$\begin{aligned} \frac{[\sigma(\langle \mathbf{x}, \mathbf{w}_j \rangle)]^p}{\|\mathbf{w}_j\|^{p-1}} &= \sigma(\langle \mathbf{x}, \mathbf{w}_j \rangle) \frac{\cos^{p-1}(\mathbf{x}, \mathbf{w}_j) \|\mathbf{w}_j\|^{p-1}}{\|\mathbf{w}_j\|^{p-1}} \\ &= \sigma(\langle \mathbf{x}, \mathbf{w}_j \rangle) \cos^{p-1}(\mathbf{x}, \mathbf{w}_j), \end{aligned}$$

that is, for each neuron  $\mathbf{w}_j$ , when compared to ReLU activation ( $p = 1$ ), the post-activation value is much smaller (penalized) if the angle separation between  $\mathbf{x}$  and  $\mathbf{w}_j$  is large. Such penalties, when  $p > 1$ , promote the alignment between training data samples and neurons, effectively altering the implicit bias of the gradient flow training, and resulting in trained networks that are more robust. We further elaborate on this point in Section 4. In addition, such generalized ReLU activation keeps the function  $f_p$  to be *2-positive-homogeneous* w.r.t. its parameters  $\boldsymbol{\theta}$ ,

i.e.,  $f_p(\mathbf{x}; \gamma \boldsymbol{\theta}) = \gamma^2 f_p(\mathbf{x}; \boldsymbol{\theta})$  for any  $\boldsymbol{\theta}$  and  $\gamma > 0$ . Many existing analyses (Du et al., 2018; Lyu et al., 2021; Kumar & Haupt, 2024) on 2-positive-homogeneous networks apply to the generalized pReLU networks.

### 3. Main Results on Adversarial Robustness

This section considers two distinct classifiers for  $\mathcal{D}_{X,Y}$ :

$$F^{(p)}(\mathbf{x}) = \sum_{k=1}^{K_1} \sigma^p(\langle \boldsymbol{\mu}_k, \mathbf{x} \rangle) - \sum_{k=K_1+1}^K \sigma^p(\langle \boldsymbol{\mu}_k, \mathbf{x} \rangle) \quad (2)$$

and

$$F(\mathbf{x}) = \sqrt{K_1} \sigma(\langle \bar{\boldsymbol{\mu}}_+, \mathbf{x} \rangle) - \sqrt{K_2} \sigma(\langle \bar{\boldsymbol{\mu}}_-, \mathbf{x} \rangle), \quad (3)$$

where  $\bar{\boldsymbol{\mu}}_+ = \frac{1}{\sqrt{K_1}} \sum_{k=1}^{K_1} \boldsymbol{\mu}_k$  and  $\bar{\boldsymbol{\mu}}_- = \frac{1}{\sqrt{K_2}} \sum_{k=K_1+1}^K \boldsymbol{\mu}_k$  are, respectively, the average direction of the positive and negative subclass centers, with  $\bar{\boldsymbol{\mu}}_+ \in \mathbb{S}^{D-1}$  and  $\bar{\boldsymbol{\mu}}_- \in \mathbb{S}^{D-1}$ .

Based on existing analyses for the implicit bias of gradient-based training, we conjecture that both classifiers can be learned by gradient flow on shallow networks with pReLU activations and proper initialization. Before studying this conjecture, we analyze these two classifiers in more detail.

#### Realizability of $F(\mathbf{x})$ and $F^{(p)}(\mathbf{x})$ by pReLU networks.

We first notice that both  $F(\mathbf{x})$  and  $F^{(p)}(\mathbf{x})$  can be realized (non-uniquely) by  $f_p(\mathbf{x}; \{\mathbf{w}_j, v_j\}_{j=1}^h)$  with some choices for the weights  $\{\mathbf{w}_j, v_j\}_{j=1}^h$ , stated as the following claim. (Verifying this claim is straightforward and we leave it to Appendix C)

**Claim.** *The following two statements are true:*

- When  $p = 1$ , let  $I_+$  and  $I_-$  be any two non-empty and disjoint subsets of  $[h]$ , let  $\boldsymbol{\lambda} := (\lambda_j) \in \mathbb{R}_{\geq 0}^h$  be any vector such that  $\sum_{j \in I_+} \lambda_j^2 = \sqrt{K_1}$  and  $\sum_{j \in I_-} \lambda_j^2 = \sqrt{K_2}$ , and let

$$\begin{aligned} \mathbf{w}_j &= \lambda_j \bar{\boldsymbol{\mu}}_+, \quad v_j = \lambda_j, \quad \forall j \in I_+, \\ \mathbf{w}_j &= \lambda_j \bar{\boldsymbol{\mu}}_-, \quad v_j = -\lambda_j, \quad \forall j \in I_-, \\ \mathbf{w}_j &= 0, \quad v_j = 0, \quad \text{otherwise}. \end{aligned}$$

Then  $f_p(\mathbf{x}; \{\mathbf{w}_j, v_j\}_{j=1}^h) \equiv F(\mathbf{x})$ .

- When  $p \geq 1$ , let  $I_1, \dots, I_K$  be any  $K$  non-empty and disjoint subsets of  $[h]$ , let  $\boldsymbol{\lambda} := (\lambda_j) \in \mathbb{R}_{\geq 0}^h$  be any vector such that  $\sum_{j \in I_k} \lambda_j^2 = 1, \forall k \in [K]$ , and let

$$\begin{aligned} \mathbf{w}_j &= \lambda_j \boldsymbol{\mu}_k, \quad v_j = \lambda_j, \quad \forall j \in I_k, k \leq K_1, \\ \mathbf{w}_j &= \lambda_j \boldsymbol{\mu}_k, \quad v_j = -\lambda_j, \quad \forall j \in I_k, k > K_1, \\ \mathbf{w}_j &= 0, \quad v_j = 0 \quad \text{otherwise}. \end{aligned}$$

Then  $f_p(\mathbf{x}; \{\mathbf{w}_j, v_j\}_{j=1}^h) \equiv F^{(p)}(\mathbf{x})$ .

**Remark 1.** As shown in the claim, there are infinitely many parameterizations of  $f_p$  (determined by the choice of subsets of  $[h]$  and  $\lambda$ ) that lead to the same function  $F(\mathbf{x})$  (or  $F^{(p)}(\mathbf{x})$ ), due to the symmetry in the network weights. That is, the resulting network represents the same function if one re-indexes (permutes) the neurons, or if one does the following rescaling on some weights  $(\mathbf{w}_j, v_j) \rightarrow (\gamma \mathbf{w}_j, v_j/\gamma)$  for some  $j \in [j]$  and some  $\gamma > 0$ .

In other words, the classifier  $F(\mathbf{x})$  corresponds to a vanilla shallow ReLU network whose non-zero neurons are either aligned with  $\bar{\mu}_+$ , the average direction of positive subclass centers, or  $\bar{\mu}_-$ , the negative counterpart. On the contrary,  $F^{(p)}(\mathbf{x})$  corresponds to a pReLU network whose non-zero neurons are aligned with one of the subclass centers. As we will see soon, both  $F(\mathbf{x})$  and  $F^{(p)}(\mathbf{x})$  achieve high prediction accuracy on clean samples from  $\mathcal{D}_{X,Y}$ , i.e., learning subclass centers is not necessary for good generalization, but  $F^{(p)}(\mathbf{x})$  is much more robust than  $F(\mathbf{x})$  against adversarial perturbation: *learning subclass centers significantly improve robustness*.

**Generalization and Robustness of  $F(\mathbf{x})$  and  $F^{(p)}(\mathbf{x})$ .** Having established that both  $F$  and  $F^{(p)}$  are realizable by a pReLU network  $f_p$ , we now study the generalization performance and robustness of  $F$  and  $F^{(p)}$ . Under the conjecture that both  $F$  and  $F^{(p)}$  can be learned by gradient flow training on a pReLU network  $f_p$ , we expect that such generalization and robustness properties to extend to  $f_p$ .

The following result states that both  $F$  and  $F^{(p)}$  achieve good generalization performance on  $\mathcal{D}_{X,Y}$  when the dimension  $D$  of the input data  $X$  is sufficiently large and the number of subclasses  $K$  is not too large.

**Proposition 1** (Generalization on clean data). *Given classifiers  $F(\mathbf{x}), F^{(p)}(\mathbf{x})$  defined in (3),(2), and a test sample  $(\mathbf{x}, y) \sim \mathcal{D}_{X,Y}$ , we have, for some constant  $C > 0$ ,*

$$\begin{aligned} \mathbb{P}(F(\mathbf{x})y > 0) &\geq 1 - 4 \exp\left(-\frac{CD}{4\alpha^2 K}\right), \\ \mathbb{P}(F^{(p)}(\mathbf{x})y > 0) &\geq 1 - 2(K+1) \exp\left(-\frac{CD}{\alpha^2 K^2}\right), \\ &\quad \forall p \geq 1. \end{aligned}$$

However, the following theorem shows a significant difference between  $F$  and  $F^{(p)}$  regarding their adversarial robustness, when the number of subclasses  $K$  is large.

**Theorem 1** ( $l_2$ -adversarial robustness). *Given classifiers  $F(\mathbf{x}), F^{(p)}(\mathbf{x})$  defined in (3),(2), and a test sample  $(\mathbf{x}, y) \sim \mathcal{D}_{X,Y}$ , the following statement is true for the same constant  $C$  in Proposition 1:*

- *(Non-robustness of  $F(\mathbf{x})$  against  $\mathcal{O}\left(\frac{1}{\sqrt{K}}\right)$ -attack)* Let

$$d_0 := \frac{\sqrt{K_1}\bar{\mu}_+ - \sqrt{K_2}\bar{\mu}_-}{\sqrt{K}} \in \mathbb{S}^{D-1}. \text{ Then for any } \rho \geq 0,$$

$$\mathbb{P}\left(F\left(\mathbf{x} - \frac{y(1+\rho)}{\sqrt{K}} d_0\right)y > 0\right) \leq 2 \exp\left(-\frac{CD\rho^2}{K\alpha^2}\right).$$

- *(Robustness of  $F^{(p)}(\mathbf{x})$  against  $\mathcal{O}(1)$ -attack)* Let  $p \geq 2$ . Then for any  $0 \leq \delta \leq \sqrt{2}$ ,

$$\begin{aligned} \mathbb{P}\left(\min_{\|\mathbf{d}\| \leq 1} F^{(p)}\left(\mathbf{x} + \frac{\sqrt{2}-\delta}{2} \mathbf{d}\right)y > 0\right) \\ \geq 1 - 2(K+1) \exp\left(-\frac{CD\delta^2}{2K^2\alpha^2}\right). \end{aligned}$$

We refer the readers to Appendix C for the proofs for Proposition 1 and Theorem 1. Our theoretical results should be understood in the high-dimensional or low-intra-class-variance regime so that  $\frac{D}{\alpha^2} \gg K^2 \log K$ . On the one hand, Theorem 1 shows that  $F(\mathbf{x})$ , a classifier that corresponds to shallow ReLU networks whose neurons are aligned with either  $\bar{\mu}_+$  or  $\bar{\mu}_-$ , is vulnerable to an adversarial attack of radius  $\mathcal{O}\left(\frac{1}{\sqrt{K}}\right)$ , which diminishes as the number of subclasses grows. On the other hand, Theorem 1 shows that  $F^{(p)}(\mathbf{x}), p \geq 2$ , a classifier that corresponds to shallow pReLU networks whose neurons are aligned with one of the subclass centers, is  $\mathcal{O}(1)$ -robust against adversarial attacks.

**Remark 2.** Complementary results (in Appendix C.4) to Theorem 1 show that  $F(\mathbf{x})$  is robust against any attack with  $l_2$ -norm slightly smaller than  $\frac{1}{\sqrt{K}}$ , and  $F^{(p)}(\mathbf{x})$  is not robust to some attack with  $l_2$ -norm slightly larger than  $\frac{\sqrt{2}}{2}$ . Therefore,  $\frac{1}{\sqrt{K}}$  (or  $\frac{\sqrt{2}}{2}$ , resp.) is the “critical” attack radius for  $F(\mathbf{x})$  (or  $F^{(p)}(\mathbf{x})$ , resp.). We conjecture this is also the case for the associated trained networks, which we verify in Section 5.1.

#### Conjecture on the outcome of gradient flow training.

We now study the conjecture that both classifiers can be learned by gradient flow on a shallow network with pReLU activations and proper initialization. Specifically, we pose the following:

**Conjecture 1.** Suppose that the intra-subclass variance  $\alpha > 0$  is sufficiently small, that one has a training dataset of sufficiently large size, and that we run gradient flow training on  $f_p(\mathbf{x}; \boldsymbol{\theta})$ ,  $\boldsymbol{\theta} = \{\mathbf{w}_j, v_j\}_{j=1}^h$  of sufficiently large width  $h$  for sufficiently long time  $T$ , starting from random initialization of the weights with a sufficiently small initialization scale. If  $p = 1$ , then we have  $\inf_{c>0} \sup_{\mathbf{x} \in \mathbb{S}^{D-1}} |cf_p(\mathbf{x}; \boldsymbol{\theta}(T)) - F(\mathbf{x})| \ll 1$ ; If  $p \in [3, \bar{p}]$  for some  $\bar{p} > 3$ , then we have  $\inf_{c>0} \sup_{\mathbf{x} \in \mathbb{S}^{D-1}} |cf_p(\mathbf{x}; \boldsymbol{\theta}(T)) - F^{(p)}(\mathbf{x})| \ll 1$ .

Our conjecture states that with proper initialization and sufficiently long training time, gradient flow finds a network

that is, up to a scaling factor  $c > 0^1$ , close to  $F(\mathbf{x})$  in  $\ell_\infty$ -distance over  $\mathbb{S}^{D-1}$ , if  $p = 1$ . That is, when training a vanilla ReLU network, the neurons learn average directions  $\bar{\mu}_+$ , and  $\bar{\mu}_-$  of subclass centers. However, when  $p \geq 3$ , gradient flow finds a network close to  $F^{(p)}(\mathbf{x})$ , i.e., the neurons learn individual subclass centers. Note that  $p$  can not be too large, as the post activation  $\frac{[\sigma(\langle \mathbf{x}, \mathbf{w}_j \rangle)]^p}{\|\mathbf{w}_j\|^{p-1}}$  converges to  $\mathbb{1}_{\cos(\mathbf{x}, \mathbf{w}_j)=1} \cdot \langle \mathbf{x}, \mathbf{w}_j \rangle$  when  $p$  grows ( $\|\mathbf{x}\| = 1$ ), effectively zeroing out post activation values almost everywhere and also the gradient, staggering gradient flow training.

Given Conjecture 1, Theorem 1 can be used to infer the robustness of the networks  $f_p(\mathbf{x}; \boldsymbol{\theta}(T))$  obtained from gradient flow training with small initialization. When training a vanilla ReLU network in this regime, we expect the trained network to be close to  $F(\mathbf{x})$ , thus non-robust against  $\mathcal{O}\left(\frac{1}{\sqrt{K}}\right)$ -attacks. When training a pReLU network with  $p \geq 3$ , we expect the trained network to be close to  $F^{(p)}(\mathbf{x})$ , which is more robust than its ReLU counterpart.

Our conjecture is based on existing work on the implicit bias of gradient flow on shallow networks with small initialization, and we carefully explain the rationale behind it in Section 4. However, formally showing such convergence results is challenging as the data distribution  $\mathcal{D}_{X,Y}$  considered here is more complicated than those for which convergence results can be successfully shown (Boursier et al., 2022; Min et al., 2024; Chistikov et al., 2023; Wang & Ma, 2023). Instead, we provide a preliminary analysis of this conjecture. Moreover, in Section 5, we provide numerical evidence that supports our conjecture.

**Comparison with Frei et al. (2023).** Frei et al. (2023) considers a similar setting as ours with only a slight difference: their data distribution differs from ours by a scaling factor of  $\sqrt{D}$  on the input data (they also allow tiny correlations between two subclass centers) and they consider shallow ReLU network with bias. They show that any trained network by gradient flow is non-robust against  $\mathcal{O}\left(\frac{1}{\sqrt{K}}\right)$  attack, which covers any initialization, but their analysis does not characterize what the classifier the trained network corresponds to. While we consider specifically small initialization, we can at least conjecture, and numerically verify, what the corresponding classifier is at the end of the training. In addition, Frei et al. (2023) show the existence of  $\mathcal{O}(1)$ -robust ReLU network if neurons are aligned with subclasses and there is some carefully chosen bias. While it already sheds light on the need for learning subclasses,

<sup>1</sup>Since  $f_p(\mathbf{x}; \boldsymbol{\theta})$  is 2-positive-homogeneous functions w.r.t. its parameter  $\boldsymbol{\theta}$ ,  $\|\boldsymbol{\theta}(t)\| \rightarrow \infty$  as training time  $t \rightarrow \infty$  (Lyu & Li, 2019; Ji & Telgarsky, 2020), the output of  $f_p(\mathbf{x}; \boldsymbol{\theta}(T))$  can never match that of  $F(\mathbf{x})$  or  $F^{(p)}(\mathbf{x})$  without a proper choice of scaling factor  $c > 0$ . However, we note, that such a scaling factor will not change the prediction sign( $f_p(\mathbf{x}; \boldsymbol{\theta}(T))$ ).

their proposed network can not be found by gradient flow. On the contrary, our proposed robust  $F^{(p)}(\mathbf{x})$  can be obtained by gradient flow.

## 4. Discussion on Gradient Flow Training

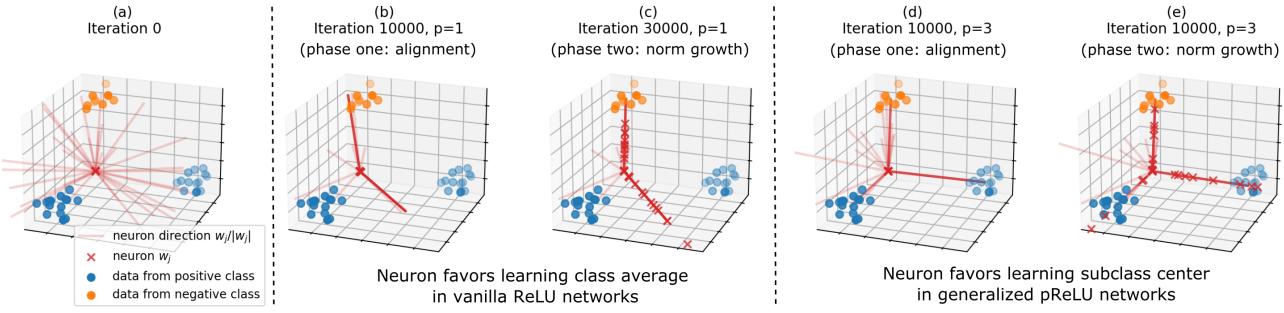
In the previous section, we conjectured that under gradient flow starting from a small initialization, a vanilla ReLU network favors learning average directions  $\bar{\mu}_+$  and  $\bar{\mu}_-$  of subclass centers, while a pReLU network favors learning every subclass centers  $\mu_k$ ,  $k \in [K]$ , resulting in a significant difference between these two networks in terms of adversarial robustness. In this section, we provide a theoretical explanation of such alignment preferences in the scope of implicit bias of gradient flow training with small initialization (Maennel et al., 2018; Boursier & Flammarion, 2024). We remind the reader that gradient flow training is described in Section 2.

### 4.1. Gradient flow with small initialization

We first briefly describe the training trajectory of gradient flow on shallow networks with a small initialization. With a sufficiently small initialization scale (to be defined later), the gradient flow training is split into two phases with distinct dynamic behaviors of the neurons. The first phase is often referred to as the initial/early *alignment phase* (Boursier & Flammarion, 2024; Min et al., 2024), during which the neurons keep small norms while changing their directions towards one of the *extremal vectors* (Maennel et al., 2018), which are jointly determined by the training dataset and the activation function. The second phase is often referred to as the fitting/convergence phase. Within the fitting phase, neurons keep a good alignment with the extremal vectors and start to grow their norms, and the loss keeps decreasing until convergence, upon which one obtains a trained network whose dominant neurons (those with large norms) are all aligned with one of the extremal vectors.

Notably, the neurons favor different extremal vectors depending on the activation function, precisely depicted by Figure 1. With the same dataset and with the same initialization of the weights, pReLU activation makes a significant difference in neuron alignment: When  $p = 1$  (vanilla ReLU network), the average class centers  $\bar{\mu}_+$  and  $\bar{\mu}_-$  are those extremal vectors “attracting” neurons during the alignment phase, leading to a trained ReLU network that has effectively two neurons (one aligned with  $\bar{\mu}_+$  and another with  $\bar{\mu}_-$ ) at the end of the training. However, when  $p = 3$ , the subclass centers  $\mu_1, \dots, \mu_k$  become extremal vectors that are “attracting” neurons, the resulting pReLU networks successfully learn every subclass center, which, we have argued in Section 3, substantially improves the robustness (over vanilla ReLU net) against adversarial attack. While the case visualized in Figure 1 is a simple one in 3-

## Can Implicit Bias Imply Adversarial Robustness?



**Figure 1.** Visualizing the training of a pReLU network under small initialization. The dataset has its positive class sampled from two subclasses. **(a)** At initialization, all neurons have small norms and point toward random directions; **When  $p = 1$** , **(b)** During the alignment phase, the neuron directions are aligned with either of the average class centers  $\bar{\mu}_+$  and  $\bar{\mu}_-$ ; **(c)** During the second phase, neurons keep the alignment with  $\bar{\mu}_+$  and  $\bar{\mu}_-$  while growing their norms; **When  $p = 3$** , **(d)** neurons learn subclass centers during alignment phase and **(e)** keep the alignment in the second phase. Note: the neurons pointing toward directions other than class/subclass centers are those not activated by any of the data points and have small norms throughout training.

dimension with 3 subclasses, we will provide experiments in higher dimension and with more subclasses in Section 5.

Despite the seemingly simple dynamic behavior of the neurons, the rigorous theoretical analysis is much more challenging due to the complicated dependence (Maennel et al., 2018; Boursier et al., 2022; Wang & Ma, 2023; Kumar & Haupt, 2024) of the extremal vectors on the dataset and the activation function. Hence prior works (Boursier et al., 2022; Min et al., 2024; Chistikov et al., 2023; Wang & Ma, 2023) assume simple training datasets and their analyses are for ReLU activation. Here, we are faced with a relatively more complex data distribution that generates our dataset, and at the same time deals with a pReLU activation, thus a full theoretical analysis of the convergence, that would prove our Conjecture 1, still has many challenges and deserves a careful treatment in a separate future work. For now, we provide some preliminary theoretical analysis of the neural alignment in pReLU networks that supports our conjecture.

### 4.2. Preliminary theoretical analysis on neuron alignment in pReLU networks

We first make the following two simplifying assumptions:

**Small and balanced initialization.** First, we obtain i.i.d. samples  $w_{j0}, j = 1, \dots, h$  drawn from some random distribution such that almost surely  $\|w_{j0}\| \leq M, \forall j$  for some  $M > 0^2$ , and then initialize the weights as

$$w_j(0) = \epsilon w_{j0}, \quad v_j(0) \in \{\|w_j(0)\|, -\|w_j(0)\|\}, \quad \forall j \in [h]. \quad (4)$$

That is,  $w_{j0}$  determines the initial direction of the neurons and we use  $\epsilon$  to control the initialization scale. This balanced assumption is standard in the analysis of shallow net-

<sup>2</sup>For example, if we sample every entries of  $w_{j0}$  by a uniform distribution over  $(-\frac{1}{\sqrt{D}}, \frac{1}{\sqrt{D}})$ .

works with small initialization (Soltanolkotabi et al., 2023; Boursier et al., 2022; Min et al., 2024).

**Simplified training dataset.** The training dataset  $\{(x_k, y_k), k = 1, \dots, K\}$  is the collection of exact subclass centers, together with their class label. That is,  $x_k = \mu_k, \quad y_k = \mathbb{1}_{\{k \leq K_1\}} - \mathbb{1}_{\{k > K_1\}}, \forall k \in [K]$ . Similar datasets with orthogonality among data points  $x_k$ s are studied in Boursier et al. (2022) for ReLU networks.

**Neuron alignment in pReLU networks.** The key to the theoretical understanding of neuron alignment is the following lemma.

**Lemma 1.** *Given some initialization from (4), if  $\epsilon = \mathcal{O}(\frac{1}{\sqrt{h}})$ , then there exists  $T = \theta(\frac{1}{K} \log \frac{1}{\sqrt{h}\epsilon})$  such that the trajectory under gradient flow training with the simplified training dataset almost surely satisfies that  $\forall t \leq T$ ,*

$$\max_j \left\| \frac{d}{dt} \frac{w_j(t)}{\|w_j(t)\|} - \text{sign}(v_j(0)) \mathcal{P}_{w_j(t)}^\perp x^{(p)}(w_j(t)) \right\| \\ = \mathcal{O}(\epsilon k \sqrt{h}),$$

where  $\mathcal{P}_{\mathbf{w}}^\perp = I - \frac{\mathbf{w}\mathbf{w}^\top}{\|\mathbf{w}\|^2}$  and

$$x^{(p)}(w) = \sum_{k=1}^K \gamma_k(w) y_k x_k \cdot p[\cos(x_k, w)]^{p-1}, \quad (5)$$

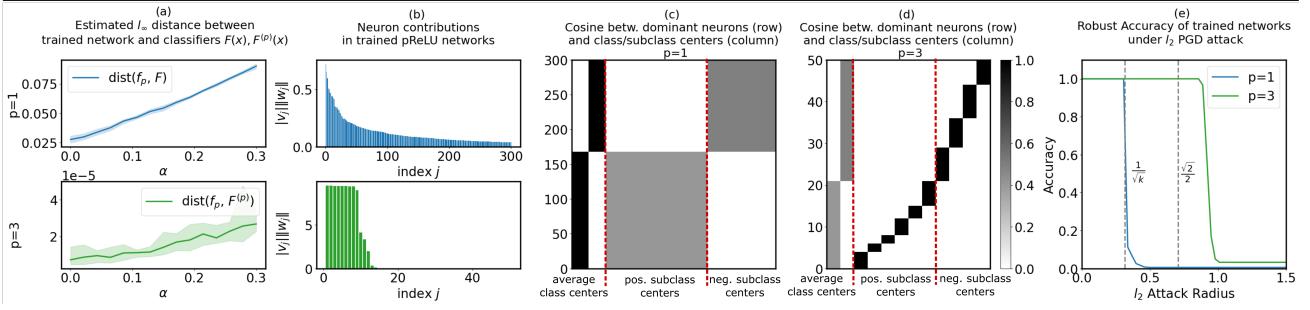
with  $\gamma_k(w)$  being a subgradient of  $\sigma(z)$  at  $z = \langle x_k, w \rangle$  when  $p = 1$  and  $\gamma_k(w) = \mathbb{1}_{\cos(x_k, w) \geq 0}$  when  $p > 1$ .

Lemma 1 suggests that during the alignment phase  $t \leq T$ , one have the following approximation

$$\frac{d}{dt} \frac{w_j(t)}{\|w_j(t)\|} \simeq \text{sign}(v_j(0)) \mathcal{P}_{w_j(t)}^\perp x^{(p)}(w_j(t)), \quad (6)$$

which essentially shows that when  $w_j$  is a *positive neuron* ( $\text{sign}(v_j(0)) > 0$ ), then gradient flow dynamics during alignment phase pushes  $w_j/\|w_j\|$  toward the direction

## Can Implicit Bias Imply Adversarial Robustness?



**Figure 2.** Numerical experiment ( $K = 10, K_1 = 6$ ) validates Conjecture 1. **(a)** We train pReLU networks using SGD with small initialization, then estimate the distance  $\text{dist}(f_p, F)$  between the trained network  $f_p$  and the classifier  $F$ , when  $p = 1$  (top plot); When  $p = 3$ , we estimate  $\text{dist}(f_p, F^{(p)})$  instead (bottom plot). The training is done under different choices of intra-subclass variance  $\alpha$  and repeated 10 times per  $\alpha$ ; the Solid line shows the average and the shade denotes the region between max and min values. **(b)** Given a trained network obtained from an instance of this training ( $\alpha = 0.1$ ), we reorder the neurons w.r.t. their contributions  $|v_j| \|\mathbf{w}_j\|$  and then plot the contributions in a bar plot; **(c)(d)** For neurons with large contributions, we plot a colormap, with each pixel represents some  $\cos(\mathbf{w}_j, \mu)$ , where  $\mu$  could be either average class centers  $\bar{\mu}_+$  and  $\bar{\mu}_-$  or subclass centers  $\mu_k, k \in [K]$ . Note: For visibility, the neurons are reordered again so that neurons aligned with the same  $\mu$  are grouped together. **(e)** Lastly, we carry out  $l_2$  PGD attack on a test dataset and plot the robust accuracy of the trained network under different choices of attack radius.

of  $\mathbf{x}^{(p)}(\mathbf{w}_j)$ . Notably,  $\mathbf{x}^{(p)}(\mathbf{w}_j)$  is a weighted combination of  $\mathbf{x}_k$ s and the mixing weights critically depend on  $p$ . Roughly speaking, when  $p = 1$ ,  $\mathbf{x}^{(p)}(\mathbf{w}_j)$ , weighing equally on  $\mathbf{x}_k$ s that activate  $\mathbf{w}_j$ , are more aligned with  $\bar{\mu}_+$  and  $\bar{\mu}_-$ , while when  $p \geq 3$ ,  $\mathbf{x}^{(p)}(\mathbf{w}_j)$ , weighing more on  $\mathbf{x}_k$ s that are close to  $\mathbf{w}_j$  in angle, are more aligned with one of the subclass centers, thus by moving toward  $\mathbf{x}^{(p)}(\mathbf{w}_j)$  in direction, the neuron  $\mathbf{w}_j$  is likely to align with average class centers in the former case, and with subclass centers in the latter case (We illustrate this in an example in Appendix D.1). This trend leads to the following alignment bias.

**Theorem 2** (Alignment bias of positive neurons). *Given some sufficiently small  $\delta > 0$  and a fixed choice of  $p \geq 1$ , then  $\exists \epsilon_0 := \epsilon_0(\delta, p) > 0$  such that for any solution of the gradient flow on  $f_p(\mathbf{x}; \theta)$  with the simplified training dataset, starting from some initialization from (4) with initialization scale  $\epsilon < \epsilon_0$ , almost surely we have that at any time  $t \leq T = \theta(\frac{1}{n} \log \frac{1}{\sqrt{h}\epsilon})$  and  $\forall j$  with  $\text{sign}(v_j(0)) > 0$ ,*

$$\frac{d}{dt} \cos(\mathbf{w}_j(t), \bar{\mu}_+) \Big|_{\cos(\mathbf{w}_j(t), \bar{\mu}_+)=1-\delta} \begin{cases} > 0, & \text{when } p = 1 \\ < 0, & \text{when } p \geq 3 \end{cases}, \quad (7)$$

and

$$\frac{d}{dt} \cos(\mathbf{w}_j(t), \mu_k) \Big|_{\cos(\mathbf{w}_j(t), \mu_k)=1-\delta} > 0, \forall k \leq K_1. \quad (8)$$

This theorem shows how different choice of  $p$  alters the neurons' preferences on which directions to align. If we define the neighborhood of some  $\mu$  direction as  $\mathcal{B}(\mu, \delta) := \{z \in \mathbb{S}^{D-1} : \cos(\mu, z) \geq 1 - \delta\}$ , then specifically for a positive neuron: When  $p = 1$ , then any neuron with  $\mathbf{w}_j(t_0) \in \mathcal{B}(\bar{\mu}_+, \delta)$  necessarily has  $\mathbf{w}_j(t) \in \mathcal{B}(\bar{\mu}_+, \delta), \forall t \in [t_0, T]$ , i.e. it keeps at least  $1 - \delta$  alignment with  $\bar{\mu}_+$  until the end of the alignment phase (at

time  $T$ ). Therefore, there is a preference of staying close to  $\bar{\mu}_+$  for positive neurons. Interestingly, such preference no longer exists when  $p \geq 3$ . In particular, any positive neuron with  $\mathbf{w}_j(t_0) \notin \mathcal{B}(\bar{\mu}_+, \delta)$  necessarily has  $\mathbf{w}_j(t) \notin \mathcal{B}(\bar{\mu}_+, \delta), \forall t \in [t_0, T]$ , i.e., any neuron initialized with some angular distance to  $\bar{\mu}_+$  will not get any closer to  $\bar{\mu}_+$ . Instead, the neurons are now more likely to stay close to some subclass centers, as shown in (8). Similar results can be said for negative neurons (whose index  $j$  has  $\text{sign}(v_j(0)) < 0$ ): they favor average negative class center  $\bar{\mu}_-$  when  $p = 1$  and favor subclass directions when  $p \geq 3$ . We refer the interested readers to Appendix D for a complete Theorem 2 (and its proof) that also includes the statement for negative neurons.

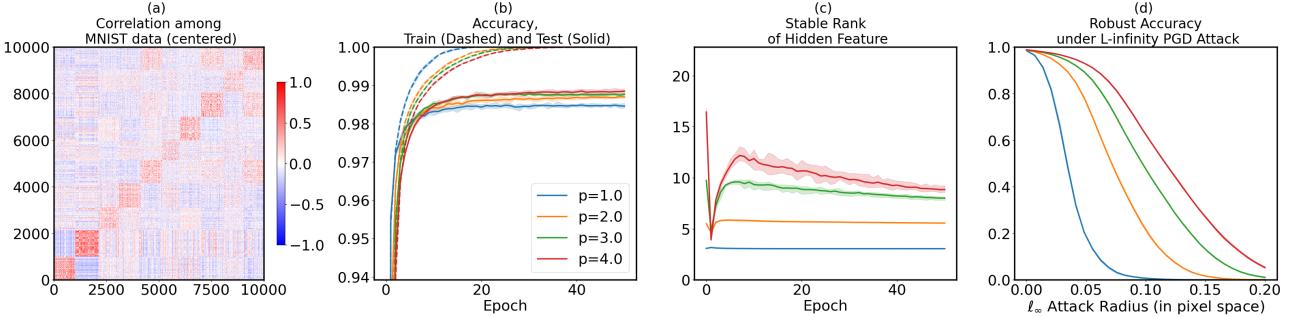
## 5. Numerical Experiments

Our numerical experiments<sup>3</sup> have two parts: First, we conduct experiments to validate Conjecture 1. Then, we provide preliminary experiments on real datasets to highlight the potential of using pReLU activation for obtaining more robust classifiers.

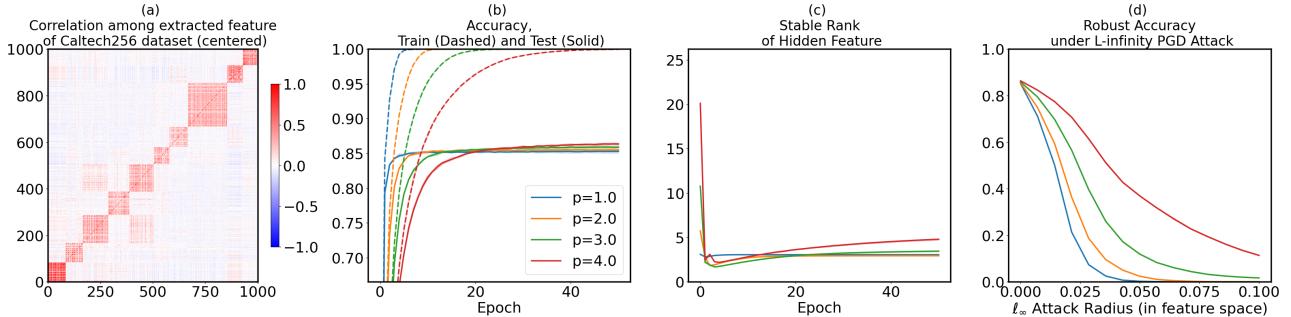
### 5.1. Numerical evidence supports our conjecture

We run SGD (batch size 100 and step size 0.2 for  $2 \times 10^5$  epochs) with small initialization (all weights initialized as mean-zero Gaussian with standard deviation  $10^{-7}$ ) to train a pReLU network with  $h = 2000$  neurons on a dataset drawn from  $\mathcal{D}_{X,Y}$  with  $D = 1000$ ,  $K = 10$ , and  $K_1 = 6$ . With different choices of intra-subclass variance  $\alpha$ , we take the network  $f_p$  at the end of the training and estimate (we refer the readers to Appendix A for the estimation algorithm) its distance  $\text{dist}(f_p, F) =$

<sup>3</sup>Code available at [https://github.com/hanchmin/pReLU\\_ICML24](https://github.com/hanchmin/pReLU_ICML24).



**Figure 3.** Parity prediction on MNIST dataset with pReLU networks. **(a)** We plot the data correlation as a colormap, where each pixel represents some  $\cos(x_i, x_j)$  between two centered data  $x_i, x_j$  from MNIST training dataset; **(b)** We run Adam with batch size 1000 to train a pReLU network under Kaiming initialization (repeated 10 times), then plot the training/testing accuracy during training for different choice of  $p$  (The shade region indicates the range between the minimum and maximum values over 10 randomized runs); **(c)** We stack the hidden post-activation representation of each training sample into a matrix and compute its stable rank, and plot the evolution of this stable rank during training; **(d)** After training for 50 epoch, we carry out APGD  $\ell_\infty$ -attack on MNIST test dataset (in pixel space) and plot the robust accuracy of the trained pReLU network under different choice of attack radius.



**Figure 4.** Classification on Caltech256 dataset (relabeled into 10 superclasses) with a pre-trained ResNet152 as a fixed feature extractor.

$\inf_{c>0} \sup_{\mathbf{x} \in \mathbb{S}^{D-1}} |cf_p(\mathbf{x}) - F(\mathbf{x})|$  to the classifier  $F(\mathbf{x})$ , or the distance  $\text{dist}(f_p, F^{(p)})$  to  $F^{(p)}(\mathbf{x})$ , depending on the value of  $p$ . As one sees from Figure 2, when  $p = 1$ ,  $f_p$  is close to  $F$ , and the estimated distance slightly increases as  $\alpha$  gets larger. Similarly, when  $p = 3$ ,  $f_p$  is close to  $F^{(p)}$ . Furthermore, we investigate the alignment between the dominant neurons in  $f_p$  and class/subclass centers. Figure 2 shows that indeed neurons in  $f_p$  learn only average class centers  $\bar{\mu}_+$  and  $\bar{\mu}_-$  when  $p = 1$  while learning every subclass center  $\mu_k, k \in [K]$  when  $p = 3$ . Lastly, as our Theorem 1 predicts,  $f_p, p = 3$  is much more robust than the one with  $p = 1$ . This series of numerical evidence strongly support Conjecture 1 (with additional experiments in Appendix A).

## 5.2. Experiments on real datasets

Although we have shown good alignment between our theory and our numerical experiments in 5.1. The settings largely remain unrealistic. To show the potential value of pReLU networks in practical settings in finding a robust classifier, we now study classification (albeit slightly modified) problems on real datasets.

### 5.2.1. PARITY CLASSIFICATION ON MNIST

We first consider training a pReLU network of width  $h = 500$  to predict whether an MNIST digit is even or odd. This poses a problem similar to the one studied in our theoretical analyses: each digit is naturally a subclass and they form classes of even digits and odd digits. Moreover, once the data is centered, as shown in Figure 3, two data points of the same digit are likely to have a large positive correlation, and two points of different digits to have a small, or negative correlation, which resembles the our orthogonality assumption among subclass centers. Therefore, with appropriate data preprocessing, we expect the pReLU network to find a more robust classifier when  $p$  is large.

**Data preprocessing.** We relabel MNIST data by parity, which leads to a binary classification. Then we center both the training and test set by subtracting off the mean image of all training data and then normalize the residue, resulting in a centered, normalized training/test set<sup>4</sup>.

**Training pReLU.** We use Kaiming initialization (He et al.,

<sup>4</sup>When training pReLU networks, some normalization of the data is required to improve training stability. To see this, notice that the post-activation value for  $i$ th data scales as  $\|x_i\|^p$ ; When  $p$  is large, this term diminishes or explodes depending on where the value  $\|x\|$  is smaller or larger than one.

2015) with non-small variance for all the weights and run Adam with cross-entropy loss and with batch size 1000 for 50 epochs and summarize the training results in Figure 3. First of all, as  $p$  increases, the trained network becomes more robust against the adversarial  $l_\infty$ -attack computed from an *adaptive projected gradient ascent* (APGD) algorithm (Croce & Hein, 2020). Interestingly, pReLU with  $p > 1$  even has a slight edge over vanilla ReLU net in terms of test accuracy on clean data. Given that the MNIST dataset does not satisfy our data distribution, there is no reason to expect that neurons in pReLU can learn each subclass (in this case, the individual digit) and indeed they cannot. However, we highlight that pReLU empirically is more capable of learning distinct patterns within each superclass/class: We stack the hidden post-activation representation of each training sample into a matrix and compute its stable rank (defined as  $\frac{\|A\|_F^2}{\|A\|^2}$  for a matrix  $A$ , as an approximation of  $\text{rank}(A)$ ). From Figure 3, we see that the hidden feature matrix of MNIST obtained from pReLU network has a much larger stable rank than the one from vanilla ReLU net, i.e. the hidden features collapse less when  $p$  is large, and we conjecture it to be the reason why pReLU still gains much more adversarial robustness than vanilla ReLU. Theoretically investigating such a phenomenon is an interesting future direction.

**Additional experiments.** To further illustrate the effectiveness of pReLU activation in improving the adversarial robustness of the trained network, we conduct additional experiments, in Appendix B, of training pReLU networks for the original digit classification task on MNIST and test the robustness of the trained network with different types of attacks. Despite that our theoretical results only suggest robustness gain can be achieved under datasets with subclasses, the additional experimental results show that the pReLU networks are still much more robust than their ReLU counterpart even for the original digit classification task, which clearly does not follow our data assumption.

#### 5.2.2. IMAGE CLASSIFICATION WITH PRE-TRAINED FEATURE EXTRACTOR

Our next experiment considers a transfer learning setting where we use some intermediate layer (more precisely, the feature layer before the fully-connected layers) of a pre-trained ResNet152 on ImageNet as a feature extractor and classify the extracted feature from Caltech256 (Griffin et al., 2007) object classification dataset with pReLU networks. The main reason behind this setting is that we expect the extracted feature representation of each class to be clustered, and this is verified in Figure 4.

We intend to have a classification task with subclasses, thus we group the original 256 classes in the dataset into 10 superclasses (each superclass has no semantic meaning in this

case) and train pReLU networks of width  $h = 2000$  that predict the superclass label. Then we obtain the feature representation of the dataset from our feature extractor, center the feature, and normalize, same as we did for the MNIST dataset. The training process is exactly the same as for the MNIST dataset, and we summarize the results in Figure 4. Even for this multi-class classification task, still pReLU achieves higher test accuracy and is more robust when  $p$  gets larger.

## 6. Conclusion

By introducing a generalized pReLU activation, we resolve the non-robustness issue, caused by the implicit bias of gradient flow, of ReLU networks when trained on a classification task in the presence of latent subclasses with small inter-subclass correlations. Future work includes formal analyses on neural alignment in pReLU networks as well as more empirical evaluation of pReLU activation for its ability to improve the robustness of neural networks.

## Acknowledgement

The authors thank the support of the NSF-Simons Research Collaborations on the Mathematical and Scientific Foundations of Deep Learning (NSF grant 2031985), and the ONR MURI Program (ONR grant 503405-78051). The authors thank the feedback from anonymous reviewers, which has greatly improved our experimental results. The authors thank Enrique Mallada, Jeremias Sulam, and Ambar Pal for their suggestions and comments at various stages of this research project, and also thank Kyle Poe for carefully reading the manuscript.

## Impact Statement

While this work is mostly theoretical, it tackles the issue of nonrobustness of neural networks trained by gradient-based algorithms, which potentially leads to more robust, reliable, and trustworthy neural networks in many application domains.

## References

- Abbe, E., Bengio, S., Boix-Adsera, E., Littwin, E., and Susskind, J. Transformers learn through gradual rank increase. *Advances in Neural Information Processing Systems*, 36, 2023.
- Athalye, A., Carlini, N., and Wagner, D. Obfuscated gradients give a false sense of security: Circumventing defenses to adversarial examples. In *International conference on machine learning*, pp. 274–283. PMLR, 2018.
- Boursier, E. and Flammarion, N. Early alignment in two-

- layer networks training is a two-edged sword. *arXiv preprint arXiv:2401.10791*, 2024.
- Boursier, E., Pullaud-Vivien, L., and Flammarion, N. 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.
- Carlini, N., Athalye, A., Papernot, N., Brendel, W., Rauber, J., Tsipras, D., Goodfellow, I., Madry, A., and Kurakin, A. On evaluating adversarial robustness. *arXiv preprint arXiv:1902.06705*, 2019.
- Chistikov, D., Englert, M., and Lazic, R. 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.
- Chizat, L. and Bach, F. Implicit bias of gradient descent for wide two-layer neural networks trained with the logistic loss. In *Proceedings of Thirty Third Conference on Learning Theory*, volume 125, pp. 1305–1338. PMLR, 09–12 Jul 2020.
- Clarke, F. H. *Optimization and nonsmooth analysis*. SIAM, 1990.
- Cohen, J., Rosenfeld, E., and Kolter, Z. Certified adversarial robustness via randomized smoothing. In *international conference on machine learning*, pp. 1310–1320. PMLR, 2019.
- Croce, F. and Hein, M. Reliable evaluation of adversarial robustness with an ensemble of diverse parameter-free attacks. In *ICML*, 2020.
- Dohmatob, E. Generalized no free lunch theorem for adversarial robustness. In *International Conference on Machine Learning*, pp. 1646–1654. PMLR, 2019.
- Du, S. S., Hu, W., and Lee, J. D. Algorithmic regularization in learning deep homogeneous models: Layers are automatically balanced. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2018.
- Faghri, F., Gowal, S., Vasconcelos, C., Fleet, D. J., Pedregosa, F., and Roux, N. L. Bridging the gap between adversarial robustness and optimization bias. *arXiv preprint arXiv:2102.08868*, 2021.
- Fawzi, A., Fawzi, H., and Fawzi, O. Adversarial vulnerability for any classifier. *Advances in neural information processing systems*, 31, 2018.
- Frei, S., Vardi, G., Bartlett, P., Srebro, N., and Hu, W. Implicit bias in leaky relu networks trained on high-dimensional data. In *The Eleventh International Conference on Learning Representations*, 2022.
- Frei, S., Vardi, G., Bartlett, P., and Srebro, N. The double-edged sword of implicit bias: Generalization vs. robustness in reLU networks. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
- Griffin, G., Holub, A., and Perona, P. Caltech-256 object category dataset. 2007.
- Gunasekar, S., Woodworth, B., Bhojanapalli, S., Neyshabur, B., and Srebro, N. Implicit regularization in matrix factorization. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, pp. 6152–6160, 2017.
- Guo, C., Rana, M., Cisse, M., and van der Maaten, L. Countering adversarial images using input transformations. In *International Conference on Learning Representations*, 2018.
- He, K., Zhang, X., Ren, S., and Sun, J. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In *Proceedings of the IEEE international conference on computer vision*, pp. 1026–1034, 2015.
- Jacot, A., Ged, F., Şimşek, B., Hongler, C., and Gabriel, F. Saddle-to-saddle dynamics in deep linear networks: Small initialization training, symmetry, and sparsity. *arXiv preprint arXiv:2106.15933*, 2021.
- Ji, Z. and Telgarsky, M. Gradient descent aligns the layers of deep linear networks. In *7th International Conference on Learning Representations, ICLR 2019*, 2019.
- Ji, Z. and Telgarsky, M. Directional convergence and alignment in deep learning. *Advances in Neural Information Processing Systems*, 33:17176–17186, 2020.
- Kumar, A. and Haupt, J. Directional convergence near small initializations and saddles in two-homogeneous neural networks. *arXiv preprint arXiv:2402.09226*, 2024.
- Levine, A. and Feizi, S. (de) randomized smoothing for certifiable defense against patch attacks. *Advances in Neural Information Processing Systems*, 33:6465–6475, 2020.
- Lyu, K. and Li, J. Gradient descent maximizes the margin of homogeneous neural networks. In *International Conference on Learning Representations*, 2019.
- Lyu, K., Li, Z., Wang, R., and Arora, S. Gradient descent on two-layer nets: Margin maximization and simplicity bias. *Advances in Neural Information Processing Systems*, 34:12978–12991, 2021.

- Maennel, H., Bousquet, O., and Gelly, S. Gradient descent quantizes relu network features. *arXiv preprint arXiv:1803.08367*, 2018.
- Min, H., Tarmoun, S., Vidal, R., and Mallada, E. On the explicit role of initialization on the convergence and implicit bias of overparametrized linear networks. In *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pp. 7760–7768. PMLR, 18–24 Jul 2021.
- Min, H., Mallada, E., and Vidal, R. Early neuron alignment in two-layer relu networks with small initialization. In *International Conference on Learning Representations*, pp. 1–8, 5 2024.
- Pal, A., Sulam, J., and Vidal, R. Adversarial examples might be avoidable: The role of data concentration in adversarial robustness. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
- Papernot, N., McDaniel, P., Wu, X., Jha, S., and Swami, A. 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.
- Salimans, T. and Kingma, D. P. Weight normalization: A simple reparameterization to accelerate training of deep neural networks. *Advances in neural information processing systems*, 29, 2016.
- Saxe, A. M., McClelland, J. L., and Ganguli, S. Exact solutions to the nonlinear dynamics of learning in deep linear neural network. In *International Conference on Learning Representations*, 2014.
- Shafahi, A., Huang, W. R., Studer, C., Feizi, S., and Goldstein, T. Are adversarial examples inevitable? In *International Conference on Learning Representations*, 2018.
- Shafahi, A., Najibi, M., Ghiasi, M. A., Xu, Z., Dickerson, J., Studer, C., Davis, L. S., Taylor, G., and Goldstein, T. Adversarial training for free! *Advances in Neural Information Processing Systems*, 32, 2019.
- Soltanolkotabi, M., Stöger, D., and Xie, C. Implicit balancing and regularization: Generalization and convergence guarantees for overparameterized asymmetric matrix sensing. In *Proceedings of Thirty Sixth Conference on Learning Theory*, volume 195, pp. 5140–5142. PMLR, 12–15 Jul 2023.
- Stöger, D. and Soltanolkotabi, M. Small random initialization is akin to spectral learning: Optimization and generalization guarantees for overparameterized low-rank matrix reconstruction. *Advances in Neural Information Processing Systems*, 34, 2021.
- Sulam, J., Muthukumar, R., and Arora, R. Adversarial robustness of supervised sparse coding. *Advances in neural information processing systems*, 33:2110–2121, 2020.
- Szegedy, C., Zaremba, W., Sutskever, I., Bruna, J., Erhan, D., Goodfellow, I., and Fergus, R. Intriguing properties of neural networks. In *2nd International Conference on Learning Representations*, 2014.
- Tarzanagh, D. A., Li, Y., Zhang, X., and Oymak, S. Max-margin token selection in attention mechanism. *Advances in Neural Information Processing Systems*, 36, 2023.
- Vardi, G., Shamir, O., and Srebro, N. On margin maximization in linear and relu networks. *Advances in Neural Information Processing Systems*, 35:37024–37036, 2022.
- Wang, M. and Ma, C. Understanding multi-phase optimization dynamics and rich nonlinear behaviors of reLU networks. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
- Wang, Z. and Jacot, A. Implicit bias of sgd in  $l_{\{2\}}$ -regularized linear dnns: One-way jumps from high to low rank. *arXiv preprint arXiv:2305.16038*, 2023.
- Wong, E., Rice, L., and Kolter, J. Z. Fast is better than free: Revisiting adversarial training. In *International Conference on Learning Representations*, 2019.
- Woodworth, B., Gunasekar, S., Lee, J. D., Moroshko, E., Savarese, P., Golan, I., Soudry, D., and Srebro, N. Kernel and rich regimes in overparametrized models. In *Conference on Learning Theory*, pp. 3635–3673. PMLR, 2020.
- Yang, G., Duan, T., Hu, J. E., Salman, H., Razenshteyn, I., and Li, J. Randomized smoothing of all shapes and sizes. In *International Conference on Machine Learning*, pp. 10693–10705. PMLR, 2020.

## A. Additional Discussion on Validating Conjecture 1

### A.1. Estimating $\text{dist}(f_p, F)$

To estimate the distance  $\text{dist}(f_p, F) = \inf_{c>0} \sup_{\mathbb{S}^{D-1}} |cf_p(\mathbf{x}) - F(\mathbf{x})|$  between a trained network  $f_p$  and a classifier  $F$  (or  $F^{(p)}$ ), we first pick an  $\hat{c} > 0$  that yields the least-square match between  $\hat{c}f_p(\mathbf{x})$  and  $F(\mathbf{x})$  over a set of points sampled from  $\text{Unif}(\mathbb{S}^{D-1})$ . Then given this choice of  $\hat{c}$ , we run normalized projected gradient ascent on  $|\hat{c}f_p(\mathbf{x}) - F(\mathbf{x})|^2$  starting from an initial  $x$  sampled from  $\text{Unif}(\mathbb{S}^{D-1})$ , and repeat a large number of times, the maximum value of  $|\hat{c}f_p(\mathbf{x}) - F(\mathbf{x})|$  at the end of Normalized PGA over all runs give an estimate of  $\text{dist}(f_p, F)$ .

---

**Algorithm 1** Estimating  $\text{dist}(f_p, F)$ 


---

**Input:** Network  $f_p$ ; Classifier  $F$  (or  $F^{(p)}$ ); step size  $\alpha$ ; sample numbers  $N_1, N_2$ ;

**Do:**

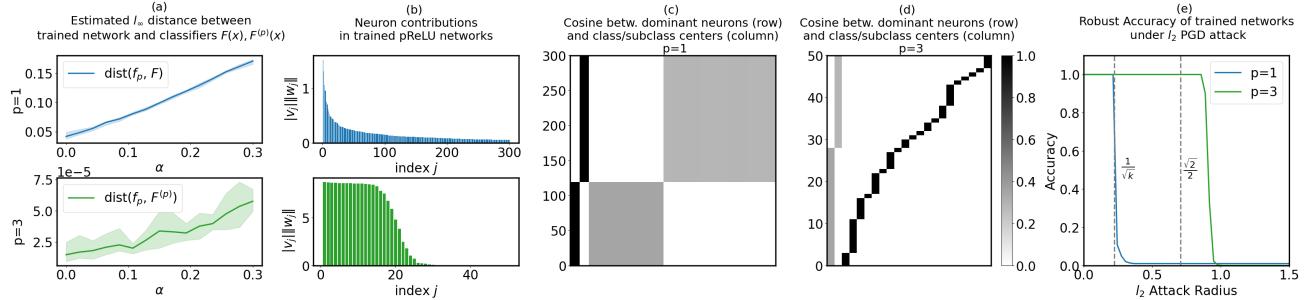
1. Sample  $\mathbf{x}, i = 1, \dots, N_1$  from  $\text{Unif}(\mathbb{S}^{D-1})$ ;  $\hat{c} \leftarrow \arg \min_{c>0} \sum_{i=1}^{N_1} |cf_p(\mathbf{x}) - F(\mathbf{x})|^2$
  2. Sample new  $\mathbf{x}, i = 1, \dots, N_2$  from  $\text{Unif}(\mathbb{S}^{D-1})$ ;  
 $\ell(\cdot) \leftarrow |\hat{c}f_p(\cdot) - F(\cdot)|^2$ ;
- For**  $i = 1, \dots, N_2$ :
- $$\mathbf{x}^{(0)} \leftarrow \mathbf{x};$$
- For**  $t = 1, \dots, \text{max\_iter}$ :
- $$\mathbf{x}^{(t)} \leftarrow \mathbf{x}^{(t-1)} + \alpha \frac{\nabla_{\mathbf{x}} \ell(\mathbf{x}^{(t-1)})}{\|\nabla_{\mathbf{x}} \ell(\mathbf{x}^{(t-1)})\|}; \quad \% \text{ normalized gradient ascent}$$
- $$\mathbf{x}^{(t)} \leftarrow \frac{\mathbf{x}^{(t)}}{\|\mathbf{x}^{(t)}\|}; \quad \% \text{ projection onto unit sphere}$$

**Return**  $\max_i |\hat{c}f_p(\mathbf{x}_i^{(\text{max\_iter})}) - F(\mathbf{x}_i^{(\text{max\_iter})})|$

---

### A.2. Additional experiment

We conduct the same experiment as in Section 5.1, with more subclasses  $K = 20, K_1 = 8$ , the results are still well aligned with Conjecture 1. Notably, as  $K$  the total number of subclasses increases, the trained ReLU network is more vulnerable to  $l_2$  attacks, compared to the case in Section 5.1, while generalized pReLU still is robust to these attacks.



**Figure 5.** Additional Numerical experiment ( $K = 20, K_1 = 8$ ) validates Conjecture 1. **(a)** We train pReLU networks using SGD with small initialization, then estimate the distance  $\text{dist}(f_p, F)$  between the trained network  $f_p$  and the classifier  $F$ , when  $p = 1$  (top plot); When  $p = 3$ , we estimate  $\text{dist}(f_p, F^{(p)})$  instead (bottom plot). The training is done under different choices of intra-subclass variance  $\alpha$  and repeated 10 times per  $\alpha$ ; the Solid line shows the average and the shade denotes the region between max and min values. **(b)** Given a trained network obtained from an instance of this training ( $\alpha = 0.1$ ), we reorder the neurons w.r.t. their contributions  $|v_j| \|w_j\|$  and then plot the contributions in a bar plot; **(c)(d)** Given neurons with large contributions, we plot a colormap, with each pixel represents some  $\cos(w_j, \mu)$ , where  $\mu$  could be either average class centers  $\bar{\mu}_+$  and  $\bar{\mu}_-$  or subclass centers  $\mu_k, k \in [K]$ . Note: For visibility, the neurons are reordered again so that neurons aligned with the same  $\mu$  are grouped together. **(e)** Lastly, we carry out  $l_2$  PGD attack on a test dataset and plot the robust accuracy of the trained network under different choices of attack radius.

## B. Additional experiments on MNIST dataset

To further illustrate the effectiveness of pReLU activation in improving the adversarial robustness of the trained network, we conduct the following additional experiments on MNIST dataset:

**MNIST digits classification.** We follow the same experiment setting in Section 5.2.1 and train the network to classify the digits instead of their parity. Figure 6 reports the training results. Despite that our theoretical results only suggest robustness gain can be achieved under datasets with subclasses, these additional experimental results show that the pReLU networks are still much more robust than their ReLU counterpart even for the original digit classification task, which clearly does not follow our data assumption.

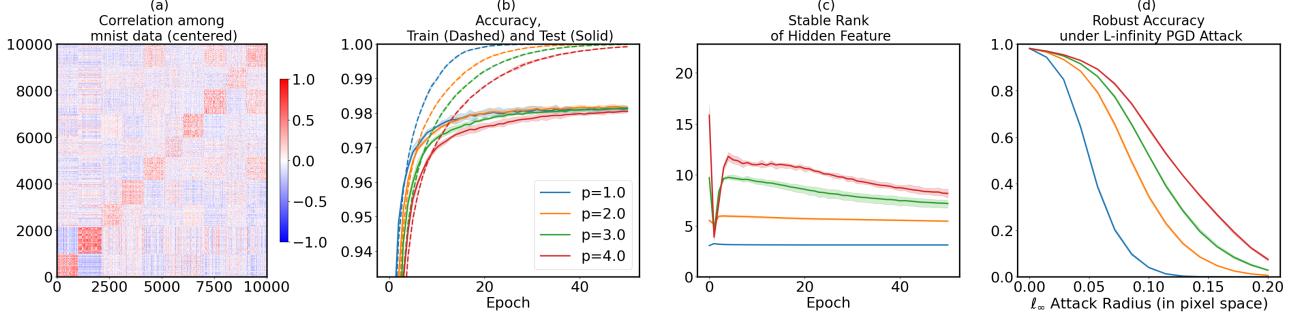


Figure 6. Digits classification on MNIST dataset with pReLU networks. (a) We plot the data correlation as a colormap, where each pixel represents some  $\cos(x_i, x_j)$  between two centered data  $x_i, x_j$  from MNIST training dataset; (b) We run Adam with batch size 1000 to train a pReLU network under Kaiming initialization (repeated 10 times), then plot the training/testing accuracy during training for different choice of  $p$ ; (c) We stack the hidden post-activation representation of each training sample into a matrix and compute its stable rank, and plot the evolution of this stable rank during training; (d) After training for 50 epoch, we carry out APGD  $l_\infty$ -attack on MNIST test dataset (in pixel space) and plot the robust accuracy of the trained pReLU network under different choice of attack radius.

**Evaluating robustness under different attacks.** Lastly, we check the adversarial robustness of the trained networks above under attacks of different norms. The results are summarized in the table below, and they show that the robustness gained from pReLU activation is agnostic to the choice of attacks.

Table 1. Robust accuracy of pReLU networks under different attacks. The reported accuracy is the mean value over 5 runs with random initialization. Bold text indicates the best accuracy within the same row.

	ReLU ( $p = 1$ )	ReLU ( $p = 2$ )	ReLU ( $p = 3$ )	ReLU ( $p = 4$ )
Clean Acc.	0.981( $\pm 0.000$ )	<b>0.982</b> ( $\pm 0.000$ )	<b>0.982</b> ( $\pm 0.000$ )	0.980 ( $\pm 0.000$ )
Robust Acc. ( $l_\infty$ , radius= 0.05)	0.512 ( $\pm 0.006$ )	0.844 ( $\pm 0.002$ )	0.892 ( $\pm 0.002$ )	<b>0.913</b> ( $\pm 0.001$ )
Robust Acc. ( $l_\infty$ , radius= 0.1)	0.040 ( $\pm 0.002$ )	0.346 ( $\pm 0.005$ )	0.522 ( $\pm 0.004$ )	<b>0.637</b> ( $\pm 0.004$ )
Robust Acc. ( $l_2$ , radius= 1)	0.301 ( $\pm 0.004$ )	0.640 ( $\pm 0.006$ )	0.726 ( $\pm 0.004$ )	<b>0.775</b> ( $\pm 0.002$ )
Robust Acc. ( $l_2$ , radius= 2)	0.007 ( $\pm 0.001$ )	0.101 ( $\pm 0.002$ )	0.165 ( $\pm 0.003$ )	<b>0.239</b> ( $\pm 0.003$ )
Robust Acc. ( $l_1$ , radius= 5)	0.500 ( $\pm 0.007$ )	0.744 ( $\pm 0.007$ )	0.780 ( $\pm 0.004$ )	<b>0.807</b> ( $\pm 0.002$ )
Robust Acc. ( $l_1$ , radius= 10)	0.098 ( $\pm 0.004$ )	0.294 ( $\pm 0.002$ )	0.354 ( $\pm 0.002$ )	<b>0.402</b> ( $\pm 0.003$ )

## C. Proofs for Proposition 1 and Theorem 1

### C.0. Verifying the claim regarding realizability of $F(\mathbf{x})$ and $F^{(p)}(\mathbf{x})$

Before we present our proofs for the main results, we start with verifying our claim regarding realizability of  $F(\mathbf{x})$  and  $F^{(p)}(\mathbf{x})$ :

**Claim** (restated). *The following two statements are true:*

- When  $p = 1$ , let  $I_+$  and  $I_-$  be any two non-empty and disjoint subsets of  $[h]$ , let  $\lambda := (\lambda_j) \in \mathbb{R}_{\geq 0}^h$  be any vector such that  $\sum_{j \in I_+} \lambda_j^2 = \sqrt{K_1}$  and  $\sum_{j \in I_-} \lambda_j^2 = \sqrt{K_2}$ , and let

$$\begin{aligned} \mathbf{w}_j &= \lambda_j \bar{\boldsymbol{\mu}}_+, v_j = \lambda_j, & \forall j \in I_+, \\ \mathbf{w}_j &= \lambda_j \bar{\boldsymbol{\mu}}_-, v_j = -\lambda_j, & \forall j \in I_-, \\ \mathbf{w}_j &= 0, v_j = 0, & \text{otherwise.} \end{aligned}$$

Then  $f_p(x; \{\mathbf{w}_j, v_j\}_{j=1}^h) \equiv F(\mathbf{x})$ .

- When  $p \geq 1$ , let  $I_1, \dots, I_K$  be any  $K$  non-empty and disjoint subsets of  $[h]$ , let  $\lambda := (\lambda_j) \in \mathbb{R}_{\geq 0}^h$  be any vector such that  $\sum_{j \in I_k} \lambda_j^2 = 1, \forall k \in [K]$ , and let

$$\begin{aligned} \mathbf{w}_j &= \lambda_j \boldsymbol{\mu}_k, v_j = \lambda_j, & \forall j \in I_k, k \leq K_1, \\ \mathbf{w}_j &= \lambda_j \boldsymbol{\mu}_k, v_j = -\lambda_j, & \forall j \in I_k, k > K_1, \\ \mathbf{w}_j &= 0, v_j = 0 & \text{otherwise.} \end{aligned}$$

Then  $f_p(x; \{\mathbf{w}_j, v_j\}_{j=1}^h) \equiv F^{(p)}(\mathbf{x})$ .

*Proof.* **Statement 1: When  $p = 1$ .** Let  $I_+, I_-, \lambda$  be defined as in the statement, then we have

$$\begin{aligned} f_p(\mathbf{x}; \{\mathbf{w}_j, v_j\}_{j=1}^h) &= \sum_{j=1}^h v_j \sigma(\langle \mathbf{x}, \mathbf{w}_j \rangle) \\ &= \sum_{j \in I_+} v_j \sigma(\langle \mathbf{x}, \mathbf{w}_j \rangle) + \sum_{j \in I_-} v_j \sigma(\langle \mathbf{x}, \mathbf{w}_j \rangle) \\ &= \sum_{j \in I_+} \lambda_j \sigma(\langle \mathbf{x}, \lambda_j \bar{\boldsymbol{\mu}}_+ \rangle) - \sum_{j \in I_-} \lambda_j \sigma(\langle \mathbf{x}, \lambda_j \bar{\boldsymbol{\mu}}_- \rangle) \\ &= \sum_{j \in I_+} \lambda_j^2 \sigma(\langle \mathbf{x}, \bar{\boldsymbol{\mu}}_+ \rangle) - \sum_{j \in I_-} \lambda_j^2 \sigma(\langle \mathbf{x}, \bar{\boldsymbol{\mu}}_- \rangle) \\ &= \sqrt{K_1} \sigma(\langle \mathbf{x}, \bar{\boldsymbol{\mu}}_+ \rangle) - \sqrt{K_2} \sigma(\langle \mathbf{x}, \bar{\boldsymbol{\mu}}_- \rangle) = F(\mathbf{x}). \end{aligned}$$

**Statement 2: When  $p \geq 1$ .** Let  $\{\mathbf{I}_k\}_{k=1}^K, \lambda$  be defined as in the statement, then we have

$$\begin{aligned} f_p(\mathbf{x}; \{\mathbf{w}_j, v_j\}_{j=1}^h) &= \sum_{j=1}^h v_j \frac{\sigma^p(\langle \mathbf{x}, \mathbf{w}_j \rangle)}{\|\mathbf{w}_j\|^{p-1}} \\ &= \sum_{k \leq K_1} \sum_{j \in I_k} v_j \frac{\sigma^p(\langle \mathbf{x}, \mathbf{w}_j \rangle)}{\|\mathbf{w}_j\|^{p-1}} + \sum_{k > K_1} \sum_{j \in I_k} v_j \frac{\sigma^p(\langle \mathbf{x}, \mathbf{w}_j \rangle)}{\|\mathbf{w}_j\|^{p-1}} \\ &= \sum_{k \leq K_1} \sum_{j \in I_k} \lambda_j \frac{\sigma^p(\langle \mathbf{x}, \lambda_j \boldsymbol{\mu}_k \rangle)}{\|\lambda_j \boldsymbol{\mu}_k\|^{p-1}} - \sum_{k > K_1} \sum_{j \in I_k} \lambda_j \frac{\sigma^p(\langle \mathbf{x}, \lambda_j \boldsymbol{\mu}_k \rangle)}{\|\lambda_j \boldsymbol{\mu}_k\|^{p-1}} \\ &= \sum_{k \leq K_1} \sum_{j \in I_k} \lambda_j^2 \sigma^p(\langle \mathbf{x}, \boldsymbol{\mu}_k \rangle) - \sum_{k > K_1} \sum_{j \in I_k} \lambda_j^2 \sigma^p(\langle \mathbf{x}, \boldsymbol{\mu}_k \rangle) \\ &= \sum_{k \leq K_1} \sigma^p(\langle \mathbf{x}, \boldsymbol{\mu}_k \rangle) - \sum_{k > K_1} \sigma^p(\langle \mathbf{x}, \boldsymbol{\mu}_k \rangle) = F^{(p)}(\mathbf{x}). \end{aligned}$$

□

### C.1. Auxiliary lemmas

The proof will use some basic facts in probability theory and we list them below.

**Lemma 2.** *Let  $\mathcal{E}_1, \mathcal{E}_2$  be two events defined on some probability space, then*

$$\mathbb{P}(\mathcal{E}_1) \leq \mathbb{P}(\mathcal{E}_1 \cap \mathcal{E}_2) + \mathbb{P}(\mathcal{E}_2^c) \quad (9)$$

*Proof.*  $\mathbb{P}(\mathcal{E}_1) = \mathbb{P}(\mathcal{E}_1 \cap (\mathcal{E}_2 \cup \mathcal{E}_2^c)) = \mathbb{P}((\mathcal{E}_1 \cap \mathcal{E}_2) \cup (\mathcal{E}_1 \cap \mathcal{E}_2^c)) \leq \mathbb{P}(\mathcal{E}_1 \cap \mathcal{E}_2) + \mathbb{P}(\mathcal{E}_1 \cap \mathcal{E}_2^c) \leq \mathbb{P}(\mathcal{E}_1 \cap \mathcal{E}_2) + \mathbb{P}(\mathcal{E}_2^c)$ .  $\square$

**Lemma 3** (Hoeffding inequality). *For any unit vector  $\mu \in \mathbb{S}^{D-1}$ , we have*

$$\mathbb{P}_{\varepsilon \sim \mathcal{N}(0, \frac{\alpha^2}{D} \mathbf{I}_D)} (|\langle \mu, \varepsilon \rangle| > t) \leq 2 \exp\left(-\frac{C D t^2}{\alpha^2}\right), \quad (10)$$

for some constant  $C > 0$ .

### C.2. Proof for Proposition 1

**Proposition 1** (Test error on clean data, restated). *Given classifiers  $F(\mathbf{x})$ ,  $F^{(p)}(\mathbf{x})$  defined in (3),(2), and a test sample  $(x, y) \sim \mathcal{D}_{X,Y}$ , we have, for some constant  $C > 0$ ,*

$$\begin{aligned} \mathbb{P}(F(\mathbf{x})y > 0) &\geq 1 - 4 \exp\left(-\frac{CD}{4\alpha^2 K}\right) \\ \mathbb{P}(F^{(p)}(\mathbf{x})y > 0) &\geq 1 - 2(K+1) \exp\left(-\frac{CD}{\alpha^2 K^2}\right), \forall p \geq 1. \end{aligned}$$

*Proof.* The proof is split into parts: We first prove the bound for  $F(\mathbf{x})$ , then show the one for  $F^{(p)}(\mathbf{x})$ . But in both cases we have

$$\mathbb{P}_{(\mathbf{x}, y) \sim \mathcal{D}_{X,Y}} (G(\mathbf{x})y > 0) = \sum_{k=1}^K \mathbb{P}_{(\mathbf{x}, y) \sim \mathcal{D}_{X,Y}} (G(\mathbf{x})y > 0 \mid z = k) \mathbb{P}(z = k), \quad (11)$$

where  $G$  can be either  $F$  or  $F^{(p)}$ . Thus, it suffices to show

$$\mathbb{P}_{(\mathbf{x}, y) \sim \mathcal{D}_{X,Y}} (G(\mathbf{x})y > 0 \mid z = k) \geq 1 - 4 \exp\left(-\frac{CD}{\alpha^2 K}\right), \forall k = 1, \dots, K. \quad (12)$$

**Bound for  $F(\mathbf{x})$**  we start with the case of  $G$  being  $F$ . When  $k \leq K_1$ , we have

$$\mathbb{P}_{(\mathbf{x}, y) \sim \mathcal{D}_{X,Y}} (F(\mathbf{x})y > 0 \mid z = k) = \mathbb{P}_{\varepsilon \sim \mathcal{N}(0, \frac{\alpha^2}{D} \mathbf{I}_D)} (F(\mu_k + \varepsilon) > 0), \quad (13)$$

and then,

$$\begin{aligned} &\mathbb{P}_{\varepsilon \sim \mathcal{N}(0, \frac{\alpha^2}{D} \mathbf{I}_D)} (F(\mu_k + \varepsilon) > 0) \\ &= 1 - \mathbb{P}_{\varepsilon \sim \mathcal{N}(0, \frac{\alpha^2}{D} \mathbf{I}_D)} (F(\mu_k + \varepsilon) < 0) \\ &= 1 - \mathbb{P}_{\varepsilon \sim \mathcal{N}(0, \frac{\alpha^2}{D} \mathbf{I}_D)} \left( \sqrt{K_1} \sigma\left(\frac{1}{\sqrt{K_1}} + \langle \varepsilon, \bar{\mu}_+ \rangle\right) - \sqrt{K_2} \sigma(\langle \varepsilon, \bar{\mu}_- \rangle) < 0, \mathcal{E}_1 \right) \\ &\geq 1 - \mathbb{P}_{\varepsilon \sim \mathcal{N}(0, \frac{\alpha^2}{D} \mathbf{I}_D)} \left( \sqrt{K_1} \left(\frac{1}{\sqrt{K_1}} - |\langle \varepsilon, \bar{\mu}_+ \rangle|\right) - \sqrt{K_2} |\langle \varepsilon, \bar{\mu}_- \rangle| < 0 \right) \end{aligned} \quad (14)$$

$$\begin{aligned} &= 1 - \mathbb{P}_{\varepsilon \sim \mathcal{N}(0, \frac{\alpha^2}{D} \mathbf{I}_D)} \left( 1 - \sqrt{K_1} |\langle \varepsilon, \bar{\mu}_+ \rangle| - \sqrt{K_2} |\langle \varepsilon, \bar{\mu}_- \rangle| < 0 \right) \\ &\geq 1 - \mathbb{P}_{\varepsilon \sim \mathcal{N}(0, \frac{\alpha^2}{D} \mathbf{I}_D)} \left( |\langle \varepsilon, \bar{\mu}_+ \rangle| > \frac{1}{2\sqrt{K_1}} \right) - \mathbb{P}_{\varepsilon \sim \mathcal{N}(0, \frac{\alpha^2}{D} \mathbf{I}_D)} \left( |\langle \varepsilon, \bar{\mu}_- \rangle| > \frac{1}{2\sqrt{K_2}} \right) \\ &\geq 1 - 2 \exp\left(-\frac{CD}{4\alpha^2 K_1}\right) - 2 \exp\left(-\frac{CD}{4\alpha^2 K_2}\right) \leq 1 - 4 \exp\left(-\frac{CD}{4\alpha^2 K}\right), \end{aligned} \quad (15)$$

where (14) uses the fact that  $\sigma(x)$  is non-decreasing w.r.t.  $x$ , and that  $\sigma(x) \leq |x|$ . The last line uses Lemma 3. The proof of the case  $k > K_1$  is identical to the one above by the symmetry of the problem.

**Bound for  $F^{(p)}(\mathbf{x})$**  When  $k \leq K_1$ , we have, again,

$$\mathbb{P}_{(\mathbf{x}, y) \sim \mathcal{D}_{X,Y}} \left( F^{(p)}(\mathbf{x})y > 0 \mid z = k \right) = \mathbb{P}_{\boldsymbol{\varepsilon} \sim \mathcal{N}\left(0, \frac{\alpha^2}{D} \mathbf{I}_D\right)} \left( F^{(p)}(\boldsymbol{\mu}_k + \boldsymbol{\varepsilon}) > 0 \right), \quad (16)$$

we define the event  $\mathcal{E}_1 := \{|\langle \boldsymbol{\mu}_k, \boldsymbol{\varepsilon} \rangle| < 1\}$ . Then by Lemma 2,

$$\begin{aligned} & \mathbb{P}_{\boldsymbol{\varepsilon} \sim \mathcal{N}\left(0, \frac{\alpha^2}{D} \mathbf{I}_D\right)} \left( F^{(p)}(\boldsymbol{\mu}_k + \boldsymbol{\varepsilon}) > 0 \right) \\ &= 1 - \mathbb{P}_{\boldsymbol{\varepsilon} \sim \mathcal{N}\left(0, \frac{\alpha^2}{D} \mathbf{I}_D\right)} \left( F^{(p)}(\boldsymbol{\mu}_k + \boldsymbol{\varepsilon}) < 0 \right) \\ &\leq 1 - \mathbb{P}_{\boldsymbol{\varepsilon} \sim \mathcal{N}\left(0, \frac{\alpha^2}{D} \mathbf{I}_D\right)} \left( F^{(p)}(\boldsymbol{\mu}_k + \boldsymbol{\varepsilon}) < 0, \mathcal{E}_1 \right) - \mathbb{P}(\mathcal{E}_1^c) \\ &\leq 1 - \mathbb{P}_{\boldsymbol{\varepsilon} \sim \mathcal{N}\left(0, \frac{\alpha^2}{D} \mathbf{I}_D\right)} \left( \sigma^p(1 + \langle \boldsymbol{\mu}_k, \boldsymbol{\varepsilon} \rangle) + \sum_{l \leq K_1, l \neq k} \sigma^p(\langle \boldsymbol{\mu}_l, \boldsymbol{\varepsilon} \rangle) - \sum_{K_1 < l \leq K_2} \sigma^p(\langle \boldsymbol{\mu}_l, \boldsymbol{\varepsilon} \rangle) < 0, \mathcal{E}_1 \right) - \mathbb{P}(\mathcal{E}_1^c) \\ &\leq 1 - \mathbb{P}_{\boldsymbol{\varepsilon} \sim \mathcal{N}\left(0, \frac{\alpha^2}{D} \mathbf{I}_D\right)} \left( (1 - |\langle \boldsymbol{\mu}_k, \boldsymbol{\varepsilon} \rangle|)^p - \sum_{l \neq k} |\langle \boldsymbol{\mu}_l, \boldsymbol{\varepsilon} \rangle|^p < 0, \mathcal{E}_1 \right) - \mathbb{P}(\mathcal{E}_1^c) \end{aligned} \quad (17)$$

$$\begin{aligned} &\leq 1 - \mathbb{P}_{\boldsymbol{\varepsilon} \sim \mathcal{N}\left(0, \frac{\alpha^2}{D} \mathbf{I}_D\right)} \left( 1 - |\langle \boldsymbol{\mu}_k, \boldsymbol{\varepsilon} \rangle| - \left( \sum_{l \neq k} |\langle \boldsymbol{\mu}_l, \boldsymbol{\varepsilon} \rangle|^p \right)^{1/p} < 0, \mathcal{E}_1 \right) - \mathbb{P}(\mathcal{E}_1^c) \\ &\leq 1 - \mathbb{P}_{\boldsymbol{\varepsilon} \sim \mathcal{N}\left(0, \frac{\alpha^2}{D} \mathbf{I}_D\right)} \left( 1 - |\langle \boldsymbol{\mu}_k, \boldsymbol{\varepsilon} \rangle| - \sum_{l \neq k} |\langle \boldsymbol{\mu}_l, \boldsymbol{\varepsilon} \rangle| < 0, \mathcal{E}_1 \right) - \mathbb{P}(\mathcal{E}_1^c) \end{aligned} \quad (18)$$

$$\begin{aligned} &\leq 1 - \mathbb{P}_{\boldsymbol{\varepsilon} \sim \mathcal{N}\left(0, \frac{\alpha^2}{D} \mathbf{I}_D\right)} \left( 1 - |\langle \boldsymbol{\mu}_k, \boldsymbol{\varepsilon} \rangle| - \sum_{l \neq k} |\langle \boldsymbol{\mu}_l, \boldsymbol{\varepsilon} \rangle| < 0 \right) - \mathbb{P}(\mathcal{E}_1^c) \\ &\leq 1 - \sum_{k=1}^K \mathbb{P}_{\boldsymbol{\varepsilon} \sim \mathcal{N}\left(0, \frac{\alpha^2}{D} \mathbf{I}_D\right)} \left( |\langle \boldsymbol{\mu}_k, \boldsymbol{\varepsilon} \rangle| > \frac{1}{K} \right) - \mathbb{P}(\mathcal{E}_1^c) \\ &\leq 1 - 2K \exp\left(-\frac{CD}{\alpha^2 K^2}\right) - 2 \exp\left(-\frac{CD}{\alpha^2}\right) \leq 1 - 2(K+1) \exp\left(-\frac{CD}{\alpha^2 K^2}\right), \end{aligned} \quad (19)$$

where (17) uses the fact that under the event  $\mathcal{E}_1$ , one has  $\sigma^p(1 + \langle \boldsymbol{\mu}_k, \boldsymbol{\varepsilon} \rangle) \geq (1 - |\langle \boldsymbol{\mu}_k, \boldsymbol{\varepsilon} \rangle|)^p$ , and (18) uses the fact that  $\|x\|_p \leq \|x\|_1$  for any vector  $x$  and  $p \geq 1$ . The last line uses Lemma 3. The proof of the case  $k > K_1$  is identical to the one above by the symmetry of the problem.  $\square$

### C.3. Proof for Theorem 1

**Theorem 1** ( $l_2$ -Adversarial Robustness, restated). *Given classifiers  $F(\mathbf{x})$ ,  $F^{(p)}(\mathbf{x})$  defined in (3),(2), and a test sample  $(x, y) \sim \mathcal{D}_{X,Y}$ , the following statement is true for some constant  $C > 0$ :*

- *(Non-robustness of  $F(\mathbf{x})$  against  $\mathcal{O}\left(\frac{1}{\sqrt{K}}\right)$ -attack)* We let  $\mathbf{d}_0 := \frac{\sqrt{K_1}\bar{\boldsymbol{\mu}}_+ - \sqrt{K_2}\bar{\boldsymbol{\mu}}_-}{\sqrt{K}} \in \mathbb{S}^{D-1}$ , then for some  $\rho > 0$ ,

$$\mathbb{P}\left(F\left(\mathbf{x} - \frac{y(1+\rho)}{\sqrt{K}} \mathbf{d}_0\right) y > 0\right) \leq 2 \exp\left(-\frac{CD\rho^2}{K\alpha^2}\right).$$

- *(Robustness of  $F^{(p)}(\mathbf{x})$  against  $\mathcal{O}(1)$ -attack)* We let  $p \geq 2$ , then for some  $0 \leq \delta < \sqrt{2}$ ,

$$\mathbb{P}\left(\min_{\|\mathbf{d}\| \leq 1} F^{(p)}\left(\mathbf{x} + \frac{\sqrt{2}-\delta}{2} \mathbf{d}\right) y > 0\right) \geq 1 - 2(K+1) \exp\left(-\frac{CD\delta^2}{2K^2\alpha^2}\right).$$

*Proof.* This proof is split into two parts. One for  $F(\mathbf{x})$  and another for  $F^{(p)}(\mathbf{x})$ .

**Non-robustness of  $F(\mathbf{x})$**  Since

$$\mathbb{P}_{(\mathbf{x}, y) \sim \mathcal{D}_{X, Y}} \left( F \left( \mathbf{x} - \frac{y(1+\rho)}{\sqrt{K}} \mathbf{d}_0 \right) y > 0 \right) = \sum_{k=1}^K \mathbb{P}_{(\mathbf{x}, y) \sim \mathcal{D}_{X, Y}} \left( F \left( \mathbf{x} - \frac{y(1+\rho)}{\sqrt{K}} \mathbf{d}_0 \right) y > 0 \mid z = k \right) \mathbb{P}(z = k), \quad (20)$$

It suffices to show

$$\mathbb{P}_{(\mathbf{x}, y) \sim \mathcal{D}_{X, Y}} \left( F \left( \mathbf{x} - \frac{y(1+\rho)}{\sqrt{K}} \mathbf{d}_0 \right) y > 0 \mid z = k \right) \leq 2 \exp \left( -\frac{CD^2 K}{(\rho-1)^2 \alpha^2} \right), \forall k \leq K \quad (21)$$

When  $k \leq K_1$ , we have

$$\begin{aligned} & \mathbb{P}_{(\mathbf{x}, y) \sim \mathcal{D}_{X, Y}} \left( F \left( \mathbf{x} - \frac{y(1+\rho)}{\sqrt{K}} \mathbf{d}_0 \right) y > 0 \mid z = k \right) \\ &= \mathbb{P}_{\boldsymbol{\varepsilon} \sim \mathcal{N}\left(0, \frac{\alpha^2}{D} \mathbf{I}_D\right)} \left( F \left( \boldsymbol{\mu}_k + \boldsymbol{\varepsilon} - \frac{1+\rho}{\sqrt{K}} \mathbf{d}_0 \right) > 0 \right) \\ &= \mathbb{P}_{\boldsymbol{\varepsilon} \sim \mathcal{N}\left(0, \frac{\alpha^2}{D} \mathbf{I}_D\right)} \left( \sqrt{K_1} \sigma \left( \frac{1}{\sqrt{K_1}} + \langle \boldsymbol{\varepsilon}, \bar{\boldsymbol{\mu}}_+ \rangle - \frac{(1+\rho)\sqrt{K_1}}{K} \right) - \sqrt{K_2} \sigma \left( \langle \boldsymbol{\varepsilon}, \bar{\boldsymbol{\mu}}_- \rangle + \frac{(1+\rho)\sqrt{K_2}}{K} \right) > 0 \right). \end{aligned} \quad (22)$$

To proceed, we define the event  $\mathcal{E}_2 := \{\langle \boldsymbol{\varepsilon}, \bar{\boldsymbol{\mu}}_+ \rangle - \frac{(1+\rho)\sqrt{K_1}}{\sqrt{K}} + \frac{1}{\sqrt{K_1}} \geq 0\}$ , then by Lemma 2

(22)

$$\begin{aligned} &\leq \mathbb{P}_{\boldsymbol{\varepsilon} \sim \mathcal{N}\left(0, \frac{\alpha^2}{D} \mathbf{I}_D\right)} \left( \sqrt{K_1} \sigma \left( \frac{1}{\sqrt{K_1}} + \langle \boldsymbol{\varepsilon}, \bar{\boldsymbol{\mu}}_+ \rangle - \frac{(1+\rho)\sqrt{K_1}}{K} \right) - \sqrt{K_2} \sigma \left( \langle \boldsymbol{\varepsilon}, \bar{\boldsymbol{\mu}}_- \rangle + \frac{(1+\rho)\sqrt{K_2}}{K} \right) > 0, \mathcal{E}_2 \right) \\ &\quad + \mathbb{P} \left( \sqrt{K_1} \sigma \left( \frac{1}{\sqrt{K_1}} + \langle \boldsymbol{\varepsilon}, \bar{\boldsymbol{\mu}}_+ \rangle - \frac{(1+\rho)\sqrt{K_1}}{K} \right) - \sqrt{K_2} \sigma \left( \langle \boldsymbol{\varepsilon}, \bar{\boldsymbol{\mu}}_- \rangle + \frac{(1+\rho)\sqrt{K_2}}{K} \right) > 0, \mathcal{E}_2^c \right) \end{aligned}$$

$$= \mathbb{P}_{\boldsymbol{\varepsilon} \sim \mathcal{N}\left(0, \frac{\alpha^2}{D} \mathbf{I}_D\right)} \left( \sqrt{K_1} \sigma \left( \frac{1}{\sqrt{K_1}} + \langle \boldsymbol{\varepsilon}, \bar{\boldsymbol{\mu}}_+ \rangle - \frac{(1+\rho)\sqrt{K_1}}{K} \right) - \sqrt{K_2} \sigma \left( \langle \boldsymbol{\varepsilon}, \bar{\boldsymbol{\mu}}_- \rangle + \frac{(1+\rho)\sqrt{K_2}}{K} \right) > 0, \mathcal{E}_2 \right) \quad (23)$$

$$\leq \mathbb{P}_{\boldsymbol{\varepsilon} \sim \mathcal{N}\left(0, \frac{\alpha^2}{D} \mathbf{I}_D\right)} \left( 1 + \sqrt{K_1} \langle \boldsymbol{\varepsilon}, \bar{\boldsymbol{\mu}}_+ \rangle - \frac{(1+\rho)K_1}{K} - \sqrt{K_2} \langle \boldsymbol{\varepsilon}, \bar{\boldsymbol{\mu}}_- \rangle - \frac{(1+\rho)K_2}{K} > 0, \mathcal{E}_2 \right) \quad (24)$$

$$\leq \mathbb{P}_{\boldsymbol{\varepsilon} \sim \mathcal{N}\left(0, \frac{\alpha^2}{D} \mathbf{I}_D\right)} \left( 1 + \sqrt{K} |\langle \boldsymbol{\varepsilon}, \mathbf{d}_0 \rangle| - (1+\rho) > 0 \right)$$

$$\leq \mathbb{P}_{\boldsymbol{\varepsilon} \sim \mathcal{N}\left(0, \frac{\alpha^2}{D} \mathbf{I}_D\right)} \left( |\langle \boldsymbol{\varepsilon}, \mathbf{d}_0 \rangle| > \frac{\rho}{\sqrt{K}} \right) \leq 2 \exp \left( -\frac{CD\rho^2}{K\alpha^2} \right), \quad (25)$$

where (23) is because the second probability is 0, and (24) uses again the monotonicity of ReLU  $\sigma(x)$ . The last line uses Lemma 3. The proof of the case  $k > K_1$  is identical to the one above by the symmetry of the problem.

**Robustness of  $F^{(p)}(\mathbf{x})$**  It suffices to show

$$\mathbb{P}_{(\mathbf{x}, y) \sim \mathcal{D}_{X, Y}} \left( \min_{\|\mathbf{d}\| \leq 1} F^{(p)} \left( \mathbf{x} + \frac{\sqrt{2}-\delta}{2} \mathbf{d} \right) y < 0 \mid z = k \right) \leq 2(K_2 + 2) \exp \left( -\frac{CD\delta^2}{2(K_2+1)^2 \alpha^2} \right). \quad (26)$$

When  $k \leq K_1$ , we have

$$\begin{aligned}
 & \mathbb{P}_{(\mathbf{x}, y) \sim \mathcal{D}_{X, Y}} \left( \min_{\|\mathbf{d}\| \leq 1} F^{(p)} \left( \mathbf{x} + \frac{\sqrt{2} - \delta}{2} \mathbf{d} \right) y < 0 \mid z = k \right) \\
 &= \mathbb{P}_{\boldsymbol{\varepsilon} \sim \mathcal{N}(0, \frac{\alpha^2}{D} \mathbf{I}_D)} \left( \min_{\|\mathbf{d}\| \leq 1} F^{(p)} \left( \boldsymbol{\mu}_k + \boldsymbol{\varepsilon} + \frac{\sqrt{2} - \delta}{2} \mathbf{d} \right) < 0 \right) \\
 &= \mathbb{P}_{\boldsymbol{\varepsilon} \sim \mathcal{N}(0, \frac{\alpha^2}{D} \mathbf{I}_D)} \left( \min_{\|\mathbf{d}\| \leq 1} \left[ \sigma^p \left( 1 + \langle \boldsymbol{\mu}_k, \boldsymbol{\varepsilon} \rangle + \frac{\sqrt{2} - \delta}{2} \langle \mathbf{d}, \boldsymbol{\mu}_k \rangle \right) \right. \right. \\
 &\quad \left. \left. + \sum_{l \neq k, 1 \leq l \leq K_1} \sigma^p \left( \langle \boldsymbol{\mu}_l, \boldsymbol{\varepsilon} \rangle + \frac{\sqrt{2} - \delta}{2} \langle \mathbf{d}, \boldsymbol{\mu}_l \rangle \right) \right] \right. \\
 &\quad \left. - \sum_{K_1+1 \leq l \leq K_2} \sigma^p \left( \langle \boldsymbol{\mu}_l, \boldsymbol{\varepsilon} \rangle + \frac{\sqrt{2} - \delta}{2} \langle \mathbf{d}, \boldsymbol{\mu}_l \rangle \right) \right] < 0 \quad (27)
 \end{aligned}$$

To proceed, we define the event

$$\mathcal{E}_3 := \left\{ 1 + \langle \boldsymbol{\mu}_k, \boldsymbol{\varepsilon} \rangle + \frac{\sqrt{2} - \delta}{2} \langle \mathbf{d}, \boldsymbol{\mu}_k \rangle \geq 0, \forall \mathbf{d} \in \mathbb{S}^{D-1} \right\}, \quad (28)$$

and for ease of presentation, we write

$$\begin{aligned}
 G(\boldsymbol{\varepsilon}, \mathbf{d}) &:= \sigma^p \left( 1 + \langle \boldsymbol{\mu}_k, \boldsymbol{\varepsilon} \rangle + \frac{\sqrt{2} - \delta}{2} \langle \mathbf{d}, \boldsymbol{\mu}_k \rangle \right) + \sum_{l \neq k, 1 \leq l \leq K_1} \sigma^p \left( \langle \boldsymbol{\mu}_l, \boldsymbol{\varepsilon} \rangle + \frac{\sqrt{2} - \delta}{2} \langle \mathbf{d}, \boldsymbol{\mu}_l \rangle \right) \\
 &\quad - \sum_{l=K_1+1}^{K_2} \sigma^p \left( \langle \boldsymbol{\mu}_l, \boldsymbol{\varepsilon} \rangle + \frac{\sqrt{2} - \delta}{2} \langle \mathbf{d}, \boldsymbol{\mu}_l \rangle \right).
 \end{aligned}$$

Then, by Lemma 2,

$$\begin{aligned}
 (27) &= \mathbb{P}_{\boldsymbol{\varepsilon} \sim \mathcal{N}(0, \frac{\alpha^2}{D} \mathbf{I}_D)} \left( \min_{\|\mathbf{d}\| \leq 1} G(\boldsymbol{\varepsilon}, \mathbf{d}) < 0 \right) \\
 &\leq \mathbb{P}_{\boldsymbol{\varepsilon} \sim \mathcal{N}(0, \frac{\alpha^2}{D} \mathbf{I}_D)} \left( \min_{\|\mathbf{d}\| \leq 1} G(\boldsymbol{\varepsilon}, \mathbf{d}) < 0, \mathcal{E}_3 \right) + \mathbb{P}(\mathcal{E}_3^c). \quad (29)
 \end{aligned}$$

Now under the event  $\mathcal{E}_3$ , we can lower bound  $G(\boldsymbol{\varepsilon}, \mathbf{d})$  by

$$\begin{aligned}
 G(\boldsymbol{\varepsilon}, \mathbf{d}) &= \sigma^p \left( 1 + \langle \boldsymbol{\mu}_k, \boldsymbol{\varepsilon} \rangle + \frac{\sqrt{2} - \delta}{2} \langle \mathbf{d}, \boldsymbol{\mu}_k \rangle \right) + \sum_{l \neq k, 1 \leq l \leq K_1} \sigma^p \left( \langle \boldsymbol{\mu}_l, \boldsymbol{\varepsilon} \rangle + \frac{\sqrt{2} - \delta}{2} \langle \mathbf{d}, \boldsymbol{\mu}_l \rangle \right) \\
 &\quad - \sum_{l=K_1+1}^{K_2} \sigma^p \left( \langle \boldsymbol{\mu}_l, \boldsymbol{\varepsilon} \rangle + \frac{\sqrt{2} - \delta}{2} \langle \mathbf{d}, \boldsymbol{\mu}_l \rangle \right). \quad (30)
 \end{aligned}$$

$$\geq \left( 1 + \langle \boldsymbol{\mu}_k, \boldsymbol{\varepsilon} \rangle + \frac{\sqrt{2} - \delta}{2} \langle \mathbf{d}, \boldsymbol{\mu}_k \rangle \right)^p - \sum_{l=K_1+1}^{K_2} \left( |\langle \boldsymbol{\mu}_l, \boldsymbol{\varepsilon} \rangle| + \frac{\sqrt{2} - \delta}{2} |\langle \mathbf{d}, \boldsymbol{\mu}_l \rangle| \right)^p. \quad (31)$$

Thus, we have

(29)

$$\begin{aligned}
 & \leq \mathbb{P}_{\boldsymbol{\varepsilon} \sim \mathcal{N}\left(0, \frac{\alpha^2}{D} \mathbf{I}_D\right)} \left( \min_{\|\mathbf{d}\| \leq 1} \left( 1 + \langle \boldsymbol{\mu}_k, \boldsymbol{\varepsilon} \rangle + \frac{\sqrt{2} - \delta}{2} \langle \mathbf{d}, \boldsymbol{\mu}_k \rangle \right)^p - \sum_{l=K_1+1}^{K_2} \left( |\langle \boldsymbol{\mu}_l, \boldsymbol{\varepsilon} \rangle| + \frac{\sqrt{2} - \delta}{2} |\langle \mathbf{d}, \boldsymbol{\mu}_l \rangle| \right)^p < 0, \mathcal{E}_3 \right) \\
 & \quad + \mathbb{P}(\mathcal{E}_3^c) \\
 & = \mathbb{P}_{\boldsymbol{\varepsilon} \sim \mathcal{N}\left(0, \frac{\alpha^2}{D} \mathbf{I}_D\right)} \left( \min_{\|\mathbf{d}\| \leq 1} 1 + \langle \boldsymbol{\mu}_k, \boldsymbol{\varepsilon} \rangle + \frac{\sqrt{2} - \delta}{2} \langle \mathbf{d}, \boldsymbol{\mu}_k \rangle - \left( \sum_{l=K_1+1}^{K_2} \left( |\langle \boldsymbol{\mu}_l, \boldsymbol{\varepsilon} \rangle| + \frac{\sqrt{2} - \delta}{2} |\langle \mathbf{d}, \boldsymbol{\mu}_l \rangle| \right)^p \right)^{\frac{1}{p}} < 0, \mathcal{E}_3 \right) \\
 & \quad + \mathbb{P}(\mathcal{E}_3^c) \\
 & \leq \mathbb{P}_{\boldsymbol{\varepsilon} \sim \mathcal{N}\left(0, \frac{\alpha^2}{D} \mathbf{I}_D\right)} \left( \underbrace{\min_{\|\mathbf{d}\| \leq 1} 1 + \frac{\sqrt{2} - \delta}{2} \langle \mathbf{d}, \boldsymbol{\mu}_k \rangle - \frac{\sqrt{2} - \delta}{2} \left( \sum_{l=K_1+1}^{K_2} (|\langle \mathbf{d}, \boldsymbol{\mu}_l \rangle|)^p \right)^{\frac{1}{p}}}_{M^*(\delta)} \right. \\
 & \quad \left. - \left( \sum_{l=K_1+1}^{K_2} (|\langle \boldsymbol{\mu}_l, \boldsymbol{\varepsilon} \rangle|)^p \right)^{\frac{1}{p}} - |\langle \boldsymbol{\mu}_k, \boldsymbol{\varepsilon} \rangle| < 0, \mathcal{E}_3 \right) + \mathbb{P}(\mathcal{E}_3^c) \tag{32}
 \end{aligned}$$

$$\begin{aligned}
 & \leq \mathbb{P}_{\boldsymbol{\varepsilon} \sim \mathcal{N}\left(0, \frac{\alpha^2}{D} \mathbf{I}_D\right)} \left( |\langle \boldsymbol{\mu}_k, \boldsymbol{\varepsilon} \rangle| + \left( \sum_{l=K_1+1}^{K_2} (|\langle \boldsymbol{\mu}_l, \boldsymbol{\varepsilon} \rangle|)^p \right)^{\frac{1}{p}} > M^*(\delta) \right) + \mathbb{P}(\mathcal{E}_3^c) \\
 & \leq \mathbb{P}_{\boldsymbol{\varepsilon} \sim \mathcal{N}\left(0, \frac{\alpha^2}{D} \mathbf{I}_D\right)} \left( |\langle \boldsymbol{\mu}_k, \boldsymbol{\varepsilon} \rangle| + \sum_{l=K_1+1}^{K_2} |\langle \boldsymbol{\mu}_l, \boldsymbol{\varepsilon} \rangle| > M^*(\delta) \right) + \mathbb{P}(\mathcal{E}_3^c) \tag{33}
 \end{aligned}$$

$$\begin{aligned}
 & \leq \mathbb{P}_{\boldsymbol{\varepsilon} \sim \mathcal{N}\left(0, \frac{\alpha^2}{D} \mathbf{I}_D\right)} \left( |\langle \boldsymbol{\mu}_k, \boldsymbol{\varepsilon} \rangle| > \frac{M^*(\delta)}{K_2 + 1} \right) + \sum_{l=K_1+1}^{K_2} \mathbb{P}_{\boldsymbol{\varepsilon} \sim \mathcal{N}\left(0, \frac{\alpha^2}{D} \mathbf{I}_D\right)} \left( |\langle \boldsymbol{\mu}_l, \boldsymbol{\varepsilon} \rangle| > \frac{M^*(\delta)}{K_2 + 1} \right) + \mathbb{P}(\mathcal{E}_3^c) \\
 & \leq \mathbb{P}_{\boldsymbol{\varepsilon} \sim \mathcal{N}\left(0, \frac{\alpha^2}{D} \mathbf{I}_D\right)} \left( |\langle \boldsymbol{\mu}_k, \boldsymbol{\varepsilon} \rangle| > \frac{M^*(\delta)}{K_2 + 1} \right) + \sum_{l=K_1+1}^{K_2} \mathbb{P}_{\boldsymbol{\varepsilon} \sim \mathcal{N}\left(0, \frac{\alpha^2}{D} \mathbf{I}_D\right)} \left( |\langle \boldsymbol{\mu}_l, \boldsymbol{\varepsilon} \rangle| > \frac{M^*(\delta)}{K_2 + 1} \right) \\
 & \quad + \mathbb{P}_{\boldsymbol{\varepsilon} \sim \mathcal{N}\left(0, \frac{\alpha^2}{D} \mathbf{I}_D\right)} \left( |\langle \boldsymbol{\mu}_k, \boldsymbol{\varepsilon} \rangle| > 1 - \frac{\sqrt{2}}{2} \right) \\
 & \leq 2(K_2 + 1) \exp \left( -\frac{CD(M^*(\delta))^2}{(K_2 + 1)^2 \alpha^2} \right) + 2 \exp \left( -\frac{CD}{4\alpha^2} \right) \leq 2(K_2 + 2) \exp \left( -\frac{CD(M^*(\delta))^2}{(K_2 + 1)^2 \alpha^2} \right), \tag{34}
 \end{aligned}$$

where (32) uses the sub-additivity of  $\ell_p$  norm, and (33) uses again the inequality  $\|\mathbf{x}\|_p \leq \|\mathbf{x}\|_1$  for any  $\mathbf{x}$  and  $p \geq 1$ . The last line uses Lemma 3. The remaining thing is to show that  $M^*(\delta) = \frac{\delta}{\sqrt{2}}$ . First, by the property of  $\ell_p$  norm (when  $p \geq 2$ ), we have

$$\left( \sum_{l=K_1+1}^{K_2} (|\langle \mathbf{d}, \boldsymbol{\mu}_l \rangle|)^p \right)^{\frac{1}{p}} \leq \sqrt{\sum_{l=K_1+1}^{K_2} |\langle \mathbf{d}, \boldsymbol{\mu}_l \rangle|^2}, \tag{35}$$

then

$$\begin{aligned}
 M^*(\delta) &= \min_{\|\mathbf{d}\| \leq 1} 1 - \frac{\sqrt{2} - \delta}{2} \langle \mathbf{d}, \boldsymbol{\mu}_k \rangle - \frac{\sqrt{2} - \delta}{2} \left( \sum_{l=K_1+1}^{K_2} (|\langle \mathbf{d}, \boldsymbol{\mu}_l \rangle|)^p \right)^{\frac{1}{p}} \\
 &\geq \min_{\|\mathbf{d}\| \leq 1} 1 - \frac{\sqrt{2} - \delta}{2} |\langle \mathbf{d}, \boldsymbol{\mu}_k \rangle| - \frac{\sqrt{2} - \delta}{2} \sqrt{\sum_{l=K_1+1}^{K_2} |\langle \mathbf{d}, \boldsymbol{\mu}_l \rangle|^2} \\
 &\geq \min_{\|\mathbf{d}\| \leq 1} 1 - \frac{\sqrt{2} - \delta}{2} \sqrt{2} \sqrt{|\langle \mathbf{d}, \boldsymbol{\mu}_k \rangle|^2 + \sum_{l=K_1+1}^{K_2} |\langle \mathbf{d}, \boldsymbol{\mu}_l \rangle|^2} \\
 &\geq \min_{\|\mathbf{d}\| \leq 1} 1 - \frac{\sqrt{2} - \delta}{2} \sqrt{2} = \frac{\delta}{\sqrt{2}}, \tag{36}
 \end{aligned}$$

$$\geq \min_{\|\mathbf{d}\| \leq 1} 1 - \frac{\sqrt{2} - \delta}{2} \sqrt{2} = \frac{\delta}{\sqrt{2}}, \tag{37}$$

where (36) uses the fact that  $\sqrt{a} + \sqrt{b} \leq \sqrt{2(a+b)}$  for any  $a, b \geq 0$ , and (37) uses the fact that  $\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_K$  are an orthonormal basis, thus  $\sqrt{\sum_{l=1}^K |\langle \mathbf{d}, \boldsymbol{\mu}_l \rangle|^2} \leq \|\mathbf{d}\| \leq 1$ .

We have proved  $M^*(\delta) \geq \frac{\delta}{\sqrt{2}}$ , and this lower bound can be attained by  $\mathbf{d}^* = \frac{-\boldsymbol{\mu}_k + \boldsymbol{\mu}_l}{\sqrt{2}}$  for any  $K_1+1 \leq l \leq K_2$ . Therefore  $M^*(\delta) = \frac{\delta}{\sqrt{2}}$ . Finally, we have

$$\begin{aligned}
 \mathbb{P}_{(\mathbf{x}, y) \sim \mathcal{D}_{X,Y}} \left( \min_{\|\mathbf{d}\| \leq 1} F^{(p)} \left( \mathbf{x} + \frac{\sqrt{2} - \delta}{2} \mathbf{d} \right) y < 0 \mid z = k \right) &\leq 2(K_2 + 2) \exp \left( -\frac{CD(M^*(\delta))^2}{(K_2 + 1)^2 \alpha^2} \right) \\
 &= 2(K_2 + 2) \exp \left( -\frac{CD\delta^2}{2(K_2 + 1)^2 \alpha^2} \right) \\
 &\leq 2(K + 1) \exp \left( -\frac{CD\delta^2}{2K^2 \alpha^2} \right). \tag{38}
 \end{aligned}$$

The proof of the case  $k > K_1$  is identical to the one above by the symmetry of the problem.  $\square$

#### C.4. Complementary results to Theorem 1

**Theorem 3** ( $l_2$ -Adversarial Robustness, complementary to Theorem 1). *Given classifiers  $F(\mathbf{x})$ ,  $F^{(p)}(\mathbf{x})$  defined in (3), (2), and a test sample  $(x, y) \sim \mathcal{D}_{X,Y}$ , the following statement is true for some constant  $C > 0$ :*

- For some  $0 \leq \rho \leq 1$ ,

$$\mathbb{P} \left( \min_{\|\mathbf{d}\| \leq 1} F \left( \mathbf{x} + \frac{(1-\rho)}{\sqrt{K}} \mathbf{d} \right) y > 0 \right) \geq 1 - 2 \exp \left( -\frac{CD\rho^2}{K\alpha^2} \right).$$

- For some  $0 < \delta$ ,

$$\mathbb{P} \left( \min_{\|\mathbf{d}\| \leq 1} F^{(p)} \left( \mathbf{x} + \frac{\sqrt{2} + \delta}{2} \mathbf{d} \right) y > 0 \right) \leq 4 \exp \left( \frac{CD\delta^2}{8K^2 \alpha^2} \right).$$

The proof has the same spirit as those for Theorem 1 so we state it briefly.

*Proof. Robustness of  $F(\mathbf{x})$ , complementary result* It suffices to show that

$$\mathbb{P}_{(\mathbf{x}, y) \sim \mathcal{D}_{X,Y}} \left( \min_{\|\mathbf{d}\| \leq 1} F \left( \mathbf{x} + \frac{(1-\rho)}{\sqrt{K}} \mathbf{d} \right) y < 0 \mid z = k \right) \leq 4 \exp \left( -\frac{CD\rho^2}{K\alpha^2} \right), \forall k \leq K. \tag{39}$$

When  $k \leq K_1$ , we have

$$\begin{aligned}
 & \mathbb{P}_{(\mathbf{x}, y) \sim \mathcal{D}_{X, Y}} \left( \min_{\|\mathbf{d}\| \leq 1} F \left( \mathbf{x} + \frac{(1-\rho)}{\sqrt{K}} \mathbf{x} \right) y < 0 \mid z = k \right) \\
 &= \mathbb{P}_{\varepsilon \sim \mathcal{N}(0, \frac{\alpha^2}{D} \mathbf{I}_D)} \left( \min_{\|\mathbf{d}\| \leq 1} F \left( \boldsymbol{\mu}_k + \varepsilon + \frac{(1-\rho)}{\sqrt{K}} \mathbf{x} \right) < 0 \right) \\
 &= \mathbb{P}_{\varepsilon \sim \mathcal{N}(0, \frac{\alpha^2}{D} \mathbf{I}_D)} \left( \min_{\|\mathbf{d}\| \leq 1} \left[ \sqrt{K_1} \sigma \left( \frac{1}{\sqrt{K_1}} + \langle \varepsilon, \bar{\boldsymbol{\mu}}_+ \rangle + \frac{1-\rho}{\sqrt{K}} \langle \mathbf{d}, \bar{\boldsymbol{\mu}}_+ \rangle \right) - \sqrt{K_2} \sigma \left( \langle \varepsilon, \bar{\boldsymbol{\mu}}_- \rangle + \frac{1-\rho}{\sqrt{K}} \langle \mathbf{d}, \bar{\boldsymbol{\mu}}_- \rangle \right) \right] < 0 \right) \tag{40}
 \end{aligned}$$

If we still let  $d_0 := \frac{\sqrt{K_1} \bar{\boldsymbol{\mu}}_+ - \sqrt{K_2} \bar{\boldsymbol{\mu}}_-}{\sqrt{K}} \in \mathbb{S}^{D-1}$ , then by the fact that  $|x| \geq \sigma(x) \geq x$  for any  $x$ , we have

$$\begin{aligned}
 & \leq \mathbb{P}_{\varepsilon \sim \mathcal{N}(0, \frac{\alpha^2}{D} \mathbf{I}_D)} \left( \min_{\|\mathbf{d}\| \leq 1} \left[ 1 + \sqrt{K_1} \langle \varepsilon, \bar{\boldsymbol{\mu}}_+ \rangle + \sqrt{K_1} \frac{1-\rho}{\sqrt{K}} \langle \mathbf{d}, \bar{\boldsymbol{\mu}}_+ \rangle - \sqrt{K_2} |\langle \varepsilon, \bar{\boldsymbol{\mu}}_- \rangle| - \sqrt{K_2} \frac{1-\rho}{\sqrt{K}} |\langle \mathbf{d}, \bar{\boldsymbol{\mu}}_- \rangle| \right] < 0 \right) \\
 & \leq \mathbb{P}_{\varepsilon \sim \mathcal{N}(0, \frac{\alpha^2}{D} \mathbf{I}_D)} \left( \min_{\|\mathbf{d}\| \leq 1} \left[ 1 - \sqrt{K_1} \frac{1-\rho}{\sqrt{K}} |\langle \mathbf{d}, \bar{\boldsymbol{\mu}}_+ \rangle| - \sqrt{K_2} \frac{1-\rho}{\sqrt{K}} |\langle \mathbf{d}, \bar{\boldsymbol{\mu}}_- \rangle| \right] - \sqrt{K_2} |\langle \varepsilon, \bar{\boldsymbol{\mu}}_- \rangle| - \sqrt{K_1} |\langle \varepsilon, \bar{\boldsymbol{\mu}}_+ \rangle| < 0 \right) \\
 & \leq \mathbb{P}_{\varepsilon \sim \mathcal{N}(0, \frac{\alpha^2}{D} \mathbf{I}_D)} \left( \min_{\|\mathbf{d}\| \leq 1} \left[ 1 - \frac{1-\rho}{\sqrt{K}} \sqrt{K_1 + K_2} \sqrt{|\langle \mathbf{d}, \bar{\boldsymbol{\mu}}_+ \rangle|^2 + |\langle \mathbf{d}, \bar{\boldsymbol{\mu}}_- \rangle|^2} \right] - \sqrt{K_2} |\langle \varepsilon, \bar{\boldsymbol{\mu}}_- \rangle| - \sqrt{K_1} |\langle \varepsilon, \bar{\boldsymbol{\mu}}_+ \rangle| < 0 \right) \tag{41}
 \end{aligned}$$

$$\leq \mathbb{P}_{\varepsilon \sim \mathcal{N}(0, \frac{\alpha^2}{D} \mathbf{I}_D)} \left( \min_{\|\mathbf{d}\| \leq 1} [1 - (1-\rho)] - \sqrt{K_2} |\langle \varepsilon, \bar{\boldsymbol{\mu}}_- \rangle| - \sqrt{K_1} |\langle \varepsilon, \bar{\boldsymbol{\mu}}_+ \rangle| < 0 \right) \tag{42}$$

$$\leq \mathbb{P}_{\varepsilon \sim \mathcal{N}(0, \frac{\alpha^2}{D} \mathbf{I}_D)} (\sqrt{K_2} |\langle \varepsilon, \bar{\boldsymbol{\mu}}_- \rangle| + \sqrt{K_1} |\langle \varepsilon, \bar{\boldsymbol{\mu}}_+ \rangle| > \rho)$$

$$\leq \mathbb{P}_{\varepsilon \sim \mathcal{N}(0, \frac{\alpha^2}{D} \mathbf{I}_D)} \left( |\langle \varepsilon, \bar{\boldsymbol{\mu}}_- \rangle| + |\langle \varepsilon, \bar{\boldsymbol{\mu}}_+ \rangle| > \frac{\rho}{\sqrt{K}} \right)$$

$$\leq 2 \mathbb{P}_{\varepsilon \sim \mathcal{N}(0, \frac{\alpha^2}{D} \mathbf{I}_D)} \left( |\langle \varepsilon, \bar{\boldsymbol{\mu}}_- \rangle| > \frac{\rho}{2\sqrt{K}} \right) \leq 4 \exp \left( -\frac{CD\rho^2}{4K\alpha^2} \right), \tag{43}$$

where (41) uses the fact that  $ab + cd \leq \sqrt{a^2 + c^2} \sqrt{b^2 + d^2}$  for any  $a, b, c, d \in \mathbb{R}$ , a simple application of Cauchy-Schwarz inequality, and (42) uses the fact that  $\bar{\boldsymbol{\mu}}_+, \bar{\boldsymbol{\mu}}_-$  are orthonormal, thus  $\sqrt{|\langle \mathbf{d}, \bar{\boldsymbol{\mu}}_+ \rangle|^2 + |\langle \mathbf{d}, \bar{\boldsymbol{\mu}}_- \rangle|^2} \leq \|\mathbf{d}\| \leq 1$ . The last line uses Lemma 3. The proof of the case  $k > K_1$  is identical to the one above by the symmetry of the problem.

**Non-robustness of  $F^{(p)}(\mathbf{x})$ , complementary result** It suffices to show that

$$\mathbb{P}_{(\mathbf{x}, y) \sim \mathcal{D}_{X, Y}} \left( \min_{\|\mathbf{d}\| \leq 1} F^{(p)} \left( \mathbf{x} + \frac{\sqrt{2} + \delta}{2} \mathbf{d} \right) y > 0 \mid z = k \right) \leq 4 \exp \left( \frac{CD\delta^2}{8K^2\alpha^2} \right), \forall k \leq K. \tag{44}$$

When  $k \leq K_1$ , we define  $\mathbf{d}_k = \frac{-\boldsymbol{\mu}_k + \boldsymbol{\mu}_K}{\sqrt{2}}$ , then

$$\begin{aligned}
 & \mathbb{P}_{(\mathbf{x}, y) \sim \mathcal{D}_{X, Y}} \left( \min_{\|\mathbf{d}\| \leq 1} F^{(p)} \left( \mathbf{x} + \frac{\sqrt{2} + \delta}{2} \mathbf{d} \right) y > 0 \mid z = k \right) \\
 & \leq \mathbb{P}_{(\mathbf{x}, y) \sim \mathcal{D}_{X, Y}} \left( F^{(p)} \left( \mathbf{x} + \frac{\sqrt{2} + \delta}{2} \mathbf{d}_k \right) y > 0 \mid z = k \right) \\
 & = \mathbb{P}_{\boldsymbol{\varepsilon} \sim \mathcal{N}(0, \frac{\alpha^2}{D} \mathbf{I}_D)} \left( F^{(p)} \left( \boldsymbol{\mu}_k + \boldsymbol{\varepsilon} + \frac{\sqrt{2} + \delta}{2} \mathbf{d}_k \right) > 0 \right) \\
 & \leq \mathbb{P}_{\boldsymbol{\varepsilon} \sim \mathcal{N}(0, \frac{\alpha^2}{D} \mathbf{I}_D)} \left( \sigma^p \left( \frac{1 - \delta/\sqrt{2}}{2} + \langle \boldsymbol{\varepsilon}, \boldsymbol{\mu}_k \rangle \right) - \sigma^p \left( \frac{1 + \delta/\sqrt{2}}{2} + \langle \boldsymbol{\varepsilon}, \boldsymbol{\mu}_K \rangle \right) + \sum_{l \neq k, l \neq K} |\langle \boldsymbol{\varepsilon}, \boldsymbol{\mu}_l \rangle|^p > 0 \right) \\
 & \leq \mathbb{P}_{\boldsymbol{\varepsilon} \sim \mathcal{N}(0, \frac{\alpha^2}{D} \mathbf{I}_D)} \left( \sigma^p \left( \frac{1 - \delta/\sqrt{2}}{2} + \langle \boldsymbol{\varepsilon}, \boldsymbol{\mu}_k \rangle \right) + \left( \sum_{l \neq k, l \neq K} |\langle \boldsymbol{\varepsilon}, \boldsymbol{\mu}_l \rangle| \right)^p > \sigma^p \left( \frac{1 + \delta/\sqrt{2}}{2} + \langle \boldsymbol{\varepsilon}, \boldsymbol{\mu}_K \rangle \right) \right) \quad (45)
 \end{aligned}$$

$$\leq \mathbb{P}_{\boldsymbol{\varepsilon} \sim \mathcal{N}(0, \frac{\alpha^2}{D} \mathbf{I}_D)} \left( \sigma^p \left( \frac{1 - \delta/\sqrt{2}}{2} + |\langle \boldsymbol{\varepsilon}, \boldsymbol{\mu}_k \rangle| \right) + \left( \sum_{l \neq k, l \neq K} |\langle \boldsymbol{\varepsilon}, \boldsymbol{\mu}_l \rangle| \right)^p > \sigma^p \left( \frac{1 + \delta/\sqrt{2}}{2} - |\langle \boldsymbol{\varepsilon}, \boldsymbol{\mu}_K \rangle| \right) \right) \quad (46)$$

$$\begin{aligned}
 & \leq \mathbb{P}_{\boldsymbol{\varepsilon} \sim \mathcal{N}(0, \frac{\alpha^2}{D} \mathbf{I}_D)} \left( \left( \sigma \left( \frac{1 - \delta/\sqrt{2}}{2} + |\langle \boldsymbol{\varepsilon}, \boldsymbol{\mu}_k \rangle| \right) + \sum_{l \neq k, l \neq K} |\langle \boldsymbol{\varepsilon}, \boldsymbol{\mu}_l \rangle| \right)^p > \sigma^p \left( \frac{1 + \delta/\sqrt{2}}{2} - |\langle \boldsymbol{\varepsilon}, \boldsymbol{\mu}_K \rangle| \right) \right) \\
 & \leq \mathbb{P}_{\boldsymbol{\varepsilon} \sim \mathcal{N}(0, \frac{\alpha^2}{D} \mathbf{I}_D)} \left( \sigma \left( \frac{1 - \delta/\sqrt{2}}{2} + |\langle \boldsymbol{\varepsilon}, \boldsymbol{\mu}_k \rangle| \right) + \sum_{l \neq k, l \neq K} |\langle \boldsymbol{\varepsilon}, \boldsymbol{\mu}_l \rangle| > \sigma \left( \frac{1 + \delta/\sqrt{2}}{2} - |\langle \boldsymbol{\varepsilon}, \boldsymbol{\mu}_K \rangle| \right) \right) \\
 & \leq \mathbb{P}_{\boldsymbol{\varepsilon} \sim \mathcal{N}(0, \frac{\alpha^2}{D} \mathbf{I}_D)} \left( \sigma \left( \frac{1 - \delta/\sqrt{2}}{2} + |\langle \boldsymbol{\varepsilon}, \boldsymbol{\mu}_k \rangle| \right) + \sum_{l \neq k, l \neq K} |\langle \boldsymbol{\varepsilon}, \boldsymbol{\mu}_l \rangle| > \frac{1 + \delta/\sqrt{2}}{2} - |\langle \boldsymbol{\varepsilon}, \boldsymbol{\mu}_K \rangle| \right), \quad (47)
 \end{aligned}$$

where (45) uses the fact that  $\sigma^p(x)$  is non-decreasing w.r.t.  $x$ , and (46) uses the fact that  $(a + b)^p \geq a^p + b^p$  for any  $a, b > 0$ . We define the event

$$\mathcal{E}_5 := \left\{ \frac{1 - \delta/\sqrt{2}}{2} + \langle \boldsymbol{\varepsilon}, \boldsymbol{\mu}_k \rangle > 0 \right\}. \quad (48)$$

Then by Lemma 2,

$$\begin{aligned}
 (47) &\leq \mathbb{P}_{\varepsilon \sim \mathcal{N}\left(0, \frac{\alpha^2}{D} \mathbf{I}_D\right)} \left( \frac{1 - \delta/\sqrt{2}}{2} + |\langle \varepsilon, \mu_k \rangle| + \sum_{l \neq k, l \neq K} |\langle \varepsilon, \mu_l \rangle| > \frac{1 + \delta/\sqrt{2}}{2} - |\langle \varepsilon, \mu_K \rangle|, \mathcal{E}_5 \right) \\
 &\quad + \mathbb{P}_{\varepsilon \sim \mathcal{N}\left(0, \frac{\alpha^2}{D} \mathbf{I}_D\right)} \left( |\langle \varepsilon, \mu_k \rangle| + \sum_{l \neq k, l \neq K} |\langle \varepsilon, \mu_l \rangle| > \frac{1 + \delta/\sqrt{2}}{2} - |\langle \varepsilon, \mu_K \rangle|, \mathcal{E}_5^c \right) \\
 &\leq \mathbb{P}_{\varepsilon \sim \mathcal{N}\left(0, \frac{\alpha^2}{D} \mathbf{I}_D\right)} \left( |\langle \varepsilon, \mu_K \rangle| + |\langle \varepsilon, \mu_k \rangle| + \sum_{l \neq k, l \neq K} |\langle \varepsilon, \mu_l \rangle| > \frac{\delta}{\sqrt{2}} \right) \\
 &\quad + \mathbb{P}_{\varepsilon \sim \mathcal{N}\left(0, \frac{\alpha^2}{D} \mathbf{I}_D\right)} \left( |\langle \varepsilon, \mu_K \rangle| + |\langle \varepsilon, \mu_k \rangle| + \sum_{l \neq k, l \neq K} |\langle \varepsilon, \mu_l \rangle| > \frac{1 + \delta/\sqrt{2}}{2} \right) \\
 &\leq 2\mathbb{P}_{\varepsilon \sim \mathcal{N}\left(0, \frac{\alpha^2}{D} \mathbf{I}_D\right)} \left( \sum_{1 \leq l \leq K} |\langle \varepsilon, \mu_l \rangle| > \frac{\delta}{2\sqrt{2}} \right) \\
 &\leq 2\mathbb{P}_{\varepsilon \sim \mathcal{N}\left(0, \frac{\alpha^2}{D} \mathbf{I}_D\right)} \left( |\langle \varepsilon, \mu_1 \rangle| > \frac{\delta}{2\sqrt{2}K} \right) \leq 4 \exp\left(\frac{CD\delta^2}{8K^2\alpha^2}\right). \tag{49}
 \end{aligned}$$

The last line uses Lemma 3. When  $k > K_1$ , the proof is similar with  $d_k := \frac{-\mu_k + \mu_1}{\sqrt{2}}$ .  $\square$

## D. Proofs for Lemma 1 and Theorem 2

### D.1. Alignment bias illustrated

As we showed in Lemma 1, during the alignment phase  $t \leq T$ , one have the following approximation

$$\frac{d}{dt} \frac{\mathbf{w}_j(t)}{\|\mathbf{w}_j(t)\|} \simeq \text{sign}(v_j(0)) \mathcal{P}_{\mathbf{w}_j(t)}^\perp x^{(p)}(\mathbf{w}_j(t)), \quad (50)$$

with

$$x^{(p)}(w) = \sum_{k=1}^K \gamma_k(w) y_k \mathbf{x}_k \cdot p[\cos(\mathbf{x}_k, w)]^{p-1}, \quad (51)$$

which essentially shows that when  $\mathbf{w}_j$  is a *positive neuron* ( $\text{sign}(v_j(0)) > 0$ ), then gradient flow dynamics during alignment phase pushes  $\mathbf{w}_j/\|\mathbf{w}_j\|$  toward the direction of  $x^{(p)}(w)$ .

Notably,  $x^{(p)}(\mathbf{w}_j)$  critically depends on  $p$ . Roughly speaking, when  $p = 1$ ,  $x^{(p)}(\mathbf{w}_j)$  are more aligned with  $\bar{\mu}_+$  and  $\bar{\mu}_-$ , while when  $p > 3$ ,  $x^{(p)}(\mathbf{w}_j)$  are more aligned with one of the subclass centers, thus by moving toward  $x^{(p)}(\mathbf{w}_j)$  in direction, the neurons are likely to align with average class centers in the former case, and with subclass centers in the latter case. We elaborate this statement here with a toy example.

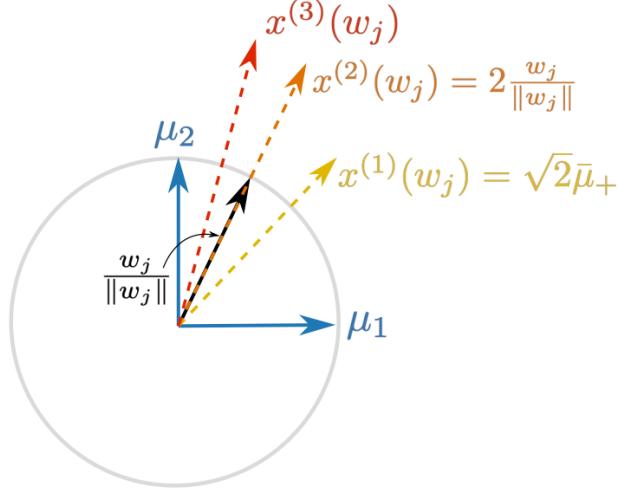


Figure 7. Alignment bias visualized. During alignment phase  $\mathbf{w}_j$  is moving toward  $x^{(p)}(\mathbf{w}_j)$  in direction. When  $p = 1$ ,  $x^{(1)}(\mathbf{w}_j)$  is aligned with average class center  $\bar{\mu}_+$ ; When  $p = 3$ ,  $x^{(p)}(\mathbf{w}_j)$  is more aligned with one of the subclass centers  $\mu_1$  and  $\mu_2$ , depending on which one is closer to  $\mathbf{w}_j$  in cosine distance.

Suppose the dataset ( $K = 2, K_1 = 2$ ) only contains two orthogonal  $\mu_1$ , and  $\mu_2$  in the 2-d plane and they both have positive labels. Given a positive neuron  $\mathbf{w}_j$  that is activated by both  $\mu_1$ , and  $\mu_2$ , as shown in Figure 7. During alignment phase  $\frac{\mathbf{w}_j}{\|\mathbf{w}_j\|}$  is moving towards the direction of  $x^{(p)}(\mathbf{w}_j)$ , which is

- when  $p = 1$ ,

$$x^{(1)}(\mathbf{w}_j) = \sum_{k=1}^K \gamma_k(\mathbf{w}_j) y_k \mathbf{x}_k = \mu_1 + \mu_2 = \sqrt{2} \bar{\mu}_+, \quad (52)$$

exactly aligned with average class center  $\bar{\mu}_+$ .

- when  $p = 2$ ,

$$x^{(2)}(\mathbf{w}_j) = \sum_{k=1}^K \gamma_k(\mathbf{w}_j) y_k \mathbf{x}_k \cdot 2[\cos(\mathbf{x}_k, w)] = 2(\mu_1 \cos(\mu_1, \mathbf{w}_j) + \mu_2 \cos(\mu_2, \mathbf{w}_j)) = 2 \frac{\mathbf{w}_j}{\|\mathbf{w}_j\|}, \quad (53)$$

exactly aligned with  $\mathbf{w}_j$  itself.

- when  $p = 3$ ,

$$x^{(3)}(w) = \sum_{k=1}^K \gamma_k(\mathbf{w}_j) y_k \mathbf{x}_k \cdot 3[\cos(\mathbf{x}_k, \mathbf{w}_j)]^2 = 3 \left( \boldsymbol{\mu}_1 \cos(\boldsymbol{\mu}_1, \mathbf{w}_j)^2 + \boldsymbol{\mu}_2 \cos(\boldsymbol{\mu}_2, \mathbf{w}_j)^2 \right), \quad (54)$$

getting closer to either  $\boldsymbol{\mu}_1$  or  $\boldsymbol{\mu}_2$ , depending which one is closer to  $\mathbf{w}_j$  in cosine distance.

Although this example is even more simplified than the one in Section 4, it is easy to visualize and keeps the core relationship between the alignment bias of the neurons and the pReLU activation. From this, we see how the alignment bias is altered under different choices of  $p$ .

## D.2. Auxiliary Lemma

We first prove the following, most analyses on gradient flow with small initialization (Boursier et al., 2022; Boursier & Flammarion, 2024; Min et al., 2024) have similar results, saying that the norm of the neurons stays close to zero during the alignment phase.

**Lemma 4.** *Given some initialization in (4), then for any  $\epsilon \leq \frac{1}{4\sqrt{h}M^2}$ , any solution to the gradient flow dynamics under the simplified training dataset satisfies*

$$\max_j \|\mathbf{w}_j(t)\|^2 \leq \frac{2\epsilon M^2}{\sqrt{h}}, \quad \max |f_p(\mathbf{x}_k; \boldsymbol{\theta}(t))| \leq 2\epsilon \sqrt{h} M^2, \quad (55)$$

$$\forall t \leq \frac{1}{4K} \log \frac{1}{\sqrt{h}\epsilon}.$$

and we need the following lemma

**Lemma 5.** *Given nonnegative  $z_1, \dots, z_n$ , consider a function*

$$g_p(q; \{z_i\}_{i=1}^n) = \left( \sum_{i=1}^n z_i^q \right) \left( \sum_{i=1}^n z_i^{p+1-q} \right), \quad (56)$$

*then  $g_p(q; \{z_i\}_{i=1}^n)$  is convex on  $\mathbb{R}$ . Moreover, as long as  $z_i \neq z_j$  for some  $i, j$ , then  $g_p(q; \{z_i\}_{i=1}^n)$  is strictly convex with minimum at  $q^* = \frac{p+1}{2}$ .*

we leave their proofs at the end of this section. Lastly, we use the following lemma from Min et al. (2024)

**Lemma 6.** *For  $\ell$  being either exponential or logistic loss, we have*

$$| -\nabla_{\hat{y}} \ell(y, \hat{y}) - y | \leq 2|\hat{y}|, \quad \forall y \in \{+1, -1\}, \quad \forall |\hat{y}| \leq 1. \quad (57)$$

## D.3. Proof for Lemma 1

**Lemma 1** (restated). *Given some initialization from (4), if  $\epsilon = \mathcal{O}(\frac{1}{\sqrt{h}})$ , then there exists  $T = \Theta(\frac{1}{K} \log \frac{1}{\sqrt{h}\epsilon})$  such that the trajectory under gradient flow training with the simplified training dataset almost surely satisfies that  $\forall t \leq T$ ,*

$$\max_j \left\| \frac{d}{dt} \frac{\mathbf{w}_j(t)}{\|\mathbf{w}_j(t)\|} - \text{sign}(v_j(0)) \mathcal{P}_{\mathbf{w}_j(t)}^\perp x^{(p)}(\mathbf{w}_j(t)) \right\| = \mathcal{O}(\epsilon k \sqrt{h}),$$

where  $\mathcal{P}_w^\perp = I - \frac{ww^\top}{\|w\|^2}$  and

$$x^{(p)}(w) = \sum_{k=1}^K \gamma_k(w) y_k \mathbf{x}_k \cdot p[\cos(\mathbf{x}_k, w)]^{p-1}, \quad (58)$$

with  $\gamma_k(w)$  being a subgradient of  $\sigma^p(z)$  at  $z = \langle \mathbf{x}_k, w \rangle$ .

*Proof.* For simplicity, we write  $\mathbf{w}_j(t)$  as  $\mathbf{w}_j$ .

As we will show in the proof for Lemma 4, under balanced initialization,

$$\frac{d}{dt} \mathbf{w}_j = - \sum_{k=1}^K \gamma_k(\mathbf{w}_j) \nabla_{\hat{y}} \ell(y_k, f_p(\mathbf{x}_k; \boldsymbol{\theta})) \|\mathbf{w}_j\| \left( \frac{p[\sigma(\langle \mathbf{x}_k, \mathbf{w}_j \rangle)]^{p-1}}{\|\mathbf{w}_j\|^{p-1}} \mathbf{x}_k - (p-1) \frac{[\sigma(\langle \mathbf{x}_k, \mathbf{w}_j \rangle)]^p}{\|\mathbf{w}_j\|^{p+1}} \mathbf{w}_j \right). \quad (59)$$

Then for any  $j \in [h]$ ,

$$\begin{aligned} \frac{d}{dt} \frac{\mathbf{w}_j}{\|\mathbf{w}_j\|} &= \mathcal{P}_{\mathbf{w}_j}^\perp \cdot \frac{1}{\|\mathbf{w}_j\|} \cdot \frac{d}{dt} \mathbf{w}_j \\ &= -\text{sign}(v_j(0)) \sum_{k=1}^K \gamma_k(\mathbf{w}_j) \nabla_{\hat{y}} \ell(y_k, f_p(\mathbf{x}_k; \boldsymbol{\theta})) \mathcal{P}_{\mathbf{w}_j}^\perp \left( p \cos(\mathbf{x}_k, \mathbf{w}_j)^{p-1} \mathbf{x}_k - (p-1) \frac{[\sigma(\langle \mathbf{x}_k, \mathbf{w}_j \rangle)]^p}{\|\mathbf{w}_j\|^{p+1}} \mathbf{w}_j \right) \\ &= -\text{sign}(v_j(0)) \sum_{i=1}^K \gamma_i(\mathbf{w}_j) \nabla_{\hat{y}} \ell(y_k, f_p(\mathbf{x}_k; \boldsymbol{\theta})) \mathcal{P}_{\mathbf{w}_j}^\perp p \cos(\mathbf{x}_k, \mathbf{w}_j)^{p-1} \mathbf{x}_k, \end{aligned}$$

Then

$$\begin{aligned} &\max_j \left\| \frac{d}{dt} \frac{\mathbf{w}_j(t)}{\|\mathbf{w}_j(t)\|} - \text{sign}(v_j(0)) \mathcal{P}_{\mathbf{w}_j(t)}^\perp x^{(p)}(\mathbf{w}_j(t)) \right\| \\ &= \max_j \left\| \sum_{i=1}^K \gamma_i(\mathbf{w}_j) \nabla_{\hat{y}} (-\ell(y_k, f_p(\mathbf{x}_k; \boldsymbol{\theta})) - y_k) \mathcal{P}_{\mathbf{w}_j}^\perp p \cos(\mathbf{x}_k, \mathbf{w}_j)^{p-1} \mathbf{x}_k \right\| \\ &\leq \max_j \left\| \sum_{i=1}^K \gamma_i(\mathbf{w}_j) |\nabla_{\hat{y}} (-\ell(y_k, f_p(\mathbf{x}_k; \boldsymbol{\theta})) - y_k)| \mathcal{P}_{\mathbf{w}_j}^\perp p \cos(\mathbf{x}_k, \mathbf{w}_j)^{p-1} \mathbf{x}_k \right\| \leq 2Kp \max_k |f_p(\mathbf{x}_k; \boldsymbol{\theta}(t))|, \end{aligned}$$

by Lemma 6. Finally, by Lemma 4, we have for any  $\epsilon \leq \frac{1}{4\sqrt{h}M^2}$ ,  $\forall t \leq T = \frac{1}{4K} \log \frac{1}{\sqrt{h}\epsilon}$ , we have

$$\max_j \left\| \frac{d}{dt} \frac{\mathbf{w}_j(t)}{\|\mathbf{w}_j(t)\|} - \text{sign}(v_j(0)) \mathcal{P}_{\mathbf{w}_j(t)}^\perp x^{(p)}(\mathbf{w}_j(t)) \right\| \leq 2Kp \max_k |f_p(\mathbf{x}_k; \boldsymbol{\theta}(t))| \leq 4\epsilon\sqrt{h}M^2Kp, \quad (60)$$

which finishes the proof.  $\square$

#### D.4. Proof for Theorem 2

**Theorem 2** (Alignment bias of neurons, complete statement). *Given some  $0 < \delta < \delta < 1$  and a fixed choice of  $p \geq 1$ , then  $\exists \epsilon_0 := \epsilon_0(\delta, p) > 0$  such that for any solution of the gradient flow on  $f_p(\mathbf{x}; \boldsymbol{\theta})$  with the simplified training dataset, starting from some initialization from (4) with initialization scale  $\epsilon < \epsilon_0$ , almost surely we have that at any time  $t \leq T = \theta(\frac{1}{n} \log \frac{1}{\sqrt{h}\epsilon})$  and*

- $\forall j$  with  $\text{sign}(v_j(0)) > 0$ ,

$$\frac{d}{dt} \cos(\mathbf{w}_j(t), \bar{\boldsymbol{\mu}}_+) \Big|_{\cos(\mathbf{w}_j(t), \bar{\boldsymbol{\mu}}_+) = 1-\delta} \begin{cases} > 0, & \text{when } p = 1 \\ < 0, & \text{when } p \geq 3 \end{cases}, \quad (61)$$

and

$$\frac{d}{dt} \cos(\mathbf{w}_j(t), \boldsymbol{\mu}_k) \Big|_{\cos(\mathbf{w}_j(t), \boldsymbol{\mu}_k) = 1-\delta} > 0, \forall k \leq K_1. \quad (62)$$

- $\forall j$  with  $\text{sign}(v_j(0)) < 0$ ,

$$\frac{d}{dt} \cos(\mathbf{w}_j(t), \bar{\boldsymbol{\mu}}_-) \Big|_{\cos(\mathbf{w}_j(t), \bar{\boldsymbol{\mu}}_-) = 1-\delta} \begin{cases} > 0, & \text{when } p = 1 \\ < 0, & \text{when } p \geq 3 \end{cases}, \quad (63)$$

and

$$\frac{d}{dt} \cos(\mathbf{w}_j(t), \boldsymbol{\mu}_k) \Big|_{\cos(\mathbf{w}_j(t), \boldsymbol{\mu}_k) = 1-\delta} > 0, \forall k > K_1. \quad (64)$$

*Proof.* The proofs for positive neurons and for negative neurons are almost identical, we will prove it for positive neurons  $\text{sign}(v_j(0)) > 0$ , i.e.  $\text{sign}(v_j(0)) = 1$ . The first part concerns about  $\frac{d}{dt} \cos(\mathbf{w}_j(t), \bar{\mu}_+) \Big|_{\cos(\mathbf{w}_j(t), \bar{\mu}_+) = 1 - \delta}$ .

$$\frac{d}{dt} \cos(\mathbf{w}_j(t), \bar{\mu}_+) \quad (65)$$

$$= \frac{1}{\|\bar{\mu}_+\|} \left\langle \frac{d}{dt} \frac{\mathbf{w}_j(t)}{\|\mathbf{w}_j(t)\|}, \bar{\mu}_+ \right\rangle \quad (66)$$

$$= \frac{1}{\|\bar{\mu}_+\|} \left\langle \mathcal{P}_{\mathbf{w}_j(t)}^\perp x^{(p)}(\mathbf{w}_j(t)), \bar{\mu}_+ \right\rangle + \frac{1}{\|\bar{\mu}_+\|} \left\langle \frac{d}{dt} \frac{\mathbf{w}_j(t)}{\|\mathbf{w}_j(t)\|} - \text{sign}(v_j(0)) \mathcal{P}_{\mathbf{w}_j(t)}^\perp x^{(p)}(\mathbf{w}_j(t)), \bar{\mu}_+ \right\rangle, \quad (67)$$

**When  $p = 1$** , if we can show that for some choice of  $\delta > 0$ ,

$$\inf_{\mathbf{w}_j \in \mathbb{S}^{D-1}, \cos(\mathbf{w}_j, \bar{\mu}_+) = 1 - \delta} \frac{1}{\|\bar{\mu}_+\|} \left\langle \mathcal{P}_{\mathbf{w}_j}^\perp x^{(1)}(\mathbf{w}_j), \bar{\mu}_+ \right\rangle := \Delta_1(\delta) > 0, \quad (68)$$

then we pick  $\epsilon \leq \epsilon_0 = \frac{\Delta_1(\delta)}{8\sqrt{h}M^2Kp}$ , and by Lemma 1, we have that for  $\forall t \leq T = \frac{1}{4K} \log \frac{1}{\sqrt{h}\epsilon}$ ,

$$\frac{d}{dt} \cos(\mathbf{w}_j(t), \bar{\mu}_+) \Big|_{\cos(\mathbf{w}_j(t), \bar{\mu}_+) = 1 - \delta} \geq \Delta_1(\delta) - \max_j \left\| \frac{d}{dt} \frac{\mathbf{w}_j(t)}{\|\mathbf{w}_j(t)\|} - \text{sign}(v_j(0)) \mathcal{P}_{\mathbf{w}_j(t)}^\perp x^{(p)}(\mathbf{w}_j(t)) \right\| \quad (69)$$

$$\geq \Delta_1(\delta) - 4\epsilon\sqrt{h}M^2Kp \geq \frac{\Delta_1(\delta)}{2} > 0, \quad (70)$$

which is what we stated in the theorem.

Similarly, **When  $p > 3$** , if we can show that for some choice of  $\delta > 0$ ,

$$\inf_{\mathbf{w}_j \in \mathbb{S}^{D-1}, \cos(\mathbf{w}_j, \bar{\mu}_+) = 1 - \delta} \frac{1}{\|\bar{\mu}_+\|} \left\langle \mathcal{P}_{\mathbf{w}_j}^\perp x^{(p)}(\mathbf{w}_j), \bar{\mu}_+ \right\rangle := \Delta_p(\delta) < 0, \quad (71)$$

then we pick  $\epsilon \leq \epsilon_0 = \frac{\Delta_p(\delta)}{8\sqrt{h}M^2Kp}$ , and by Lemma 1, we have that for  $\forall t \leq T = \frac{1}{4K} \log \frac{1}{\sqrt{h}\epsilon}$ ,

$$\frac{d}{dt} \cos(\mathbf{w}_j(t), \bar{\mu}_+) \Big|_{\cos(\mathbf{w}_j(t), \bar{\mu}_+) = 1 - \delta} \leq \Delta_p(\delta) + \max_j \left\| \frac{d}{dt} \frac{\mathbf{w}_j(t)}{\|\mathbf{w}_j(t)\|} - \text{sign}(v_j(0)) \mathcal{P}_{\mathbf{w}_j(t)}^\perp x^{(p)}(\mathbf{w}_j(t)) \right\| \quad (72)$$

$$\leq \Delta_p(\delta) + 4\epsilon\sqrt{h}M^2Kp \leq \frac{\Delta_p(\delta)}{2} < 0. \quad (73)$$

Therefore, for the first part, it suffices to show

$$\inf_{\mathbf{w}_j \in \mathbb{S}^{D-1}, \cos(\mathbf{w}_j, \bar{\mu}_+) = 1 - \delta} \frac{1}{\|\bar{\mu}_+\|} \left\langle \mathcal{P}_{\mathbf{w}_j}^\perp x^{(p)}(\mathbf{w}_j), \bar{\mu}_+ \right\rangle := \Delta_p(\delta) \begin{cases} > 0, & \text{for } p = 1 \\ > 0, & \text{for } p \geq 3 \end{cases} \quad (74)$$

Now there exists  $1 > \bar{\delta}_1 > 0$  such that when  $\delta > \bar{\delta}_1$ , and  $\cos(\mathbf{w}_j, \bar{\mu}_+) = \sqrt{1 - \delta}$ , we have  $\gamma_k(\mathbf{w}_j) = 1, \forall k \leq K_1$  and  $\gamma_k(\mathbf{w}_j) = 0, \forall k > K_1$ , i.e.,  $\mathbf{w}_j$  is activated by all  $\mathbf{x}_k, k \geq K_1$  with positive label and is not activated by any of the  $\mathbf{x}_k, k > K_1$  with negative label. Moreover, there exists  $z_k, k \leq K_1$ , such that  $\mathbf{w}_j = \|\mathbf{w}_j\| \sum_{k \leq K_1} z_k \mathbf{x}_k$  and  $z_k = \left\langle \mathbf{x}_k, \frac{\mathbf{w}_j}{\|\mathbf{w}_j\|} \right\rangle$ , i.e.,  $\mathbf{w}_j$  lies completely within the span of  $\mathbf{x}_k, k \leq K_1$ .

With this, we have

$$\begin{aligned}
 & \frac{1}{\|\bar{\mu}_+\|} \left\langle \mathcal{P}_{\mathbf{w}_j}^\perp x^{(p)}(\mathbf{w}_j), \bar{\mu}_+ \right\rangle \\
 &= \frac{1}{\|\bar{\mu}_+\|} \left\langle \left( I - \frac{\mathbf{w}_j \mathbf{w}_j^T}{\|\mathbf{w}_j\|^2} \right) \sum_k^K \gamma_k(\mathbf{w}_j) y_k \mathbf{x}_k \cdot p[\cos(\mathbf{x}_k, \mathbf{w}_j)]^{p-1}, \bar{\mu}_+ \right\rangle \\
 &= \frac{1}{K} \left\langle \left( I - \frac{\mathbf{w}_j \mathbf{w}_j^T}{\|\mathbf{w}_j\|^2} \right) \sum_{k \leq K_1} \mathbf{x}_k \cdot p[\cos(\mathbf{x}_k, \mathbf{w}_j)]^{p-1}, \sum_{k \leq K_1} \mathbf{x}_k \right\rangle \\
 &= \left\langle \sum_{k \leq K_1} \mathbf{x}_k, \sum_{k \leq K_1} \mathbf{x}_k \cdot p[\cos(\mathbf{x}_k, \mathbf{w}_j)]^{p-1} \right\rangle - \left\langle \sum_{k \leq K_1} \mathbf{x}_k, \frac{\mathbf{w}_j}{\|\mathbf{w}_j\|} \right\rangle \sum_{k \leq K_1} \left\langle \mathbf{x}_k, \frac{\mathbf{w}_j}{\|\mathbf{w}_j\|} \right\rangle \cdot p[\cos(\mathbf{x}_k, \mathbf{w}_j)]^{p-1} \\
 &= \left\langle \sum_{k \leq K_1} \mathbf{x}_k, \sum_{k \leq K_1} \mathbf{x}_k \cdot p \left\langle \mathbf{x}_k, \frac{\mathbf{w}_j}{\|\mathbf{w}_j\|} \right\rangle^{p-1} \right\rangle - \left\langle \sum_{k \leq K_1} \mathbf{x}_k, \frac{\mathbf{w}_j}{\|\mathbf{w}_j\|} \right\rangle \sum_{k \leq K_1} \left\langle \mathbf{x}_k, \frac{\mathbf{w}_j}{\|\mathbf{w}_j\|} \right\rangle \cdot p \left\langle \mathbf{x}_k, \frac{\mathbf{w}_j}{\|\mathbf{w}_j\|} \right\rangle^{p-1} \\
 &= \sum_{k \leq K_1} p \left\langle \mathbf{x}_k, \frac{\mathbf{w}_j}{\|\mathbf{w}_j\|} \right\rangle^{p-1} - \left( \sum_{k \leq K_1} \left\langle \mathbf{x}_k, \frac{\mathbf{w}_j}{\|\mathbf{w}_j\|} \right\rangle \right) \left( p \sum_{k \leq K_1} \left\langle \mathbf{x}_k, \frac{\mathbf{w}_j}{\|\mathbf{w}_j\|} \right\rangle^p \right) \\
 &= p \sum_{k \leq K_1} \left\langle \mathbf{x}_k, \frac{\mathbf{w}_j}{\|\mathbf{w}_j\|} \right\rangle^{p-1} - p \left( \sum_{k \leq K_1} \left\langle \mathbf{x}_k, \frac{\mathbf{w}_j}{\|\mathbf{w}_j\|} \right\rangle \right) \left( \sum_{k \leq K_1} \left\langle \mathbf{x}_k, \frac{\mathbf{w}_j}{\|\mathbf{w}_j\|} \right\rangle^p \right) \\
 &= p \left( \sum_{k \leq K_1} z_k^{p-1} - \left( \sum_{k \leq K_1} z_k \right) \left( \sum_{k \leq K_1} z_k^p \right) \right),
 \end{aligned}$$

Since  $\mathbf{w}_j$  lies completely within the span of  $\mathbf{x}_k, k \leq K_1$ , we have  $\sum_{k \leq K_1} z_k^2 = 1$ , then

$$\begin{aligned}
 & p \left( \sum_{k \leq K_1} z_k^{p-1} - \left( \sum_{k \leq K_1} z_k \right) \left( \sum_{k \leq K_1} z_k^p \right) \right) \\
 &= p \left( \left( \sum_{k \leq K_1} z_k^{p-1} \right) \left( \sum_{k \leq K_1} z_k^2 \right) - \left( \sum_{k \leq K_1} z_k \right) \left( \sum_{k \leq K_1} z_k^p \right) \right) = p [g_p(2; \{z_k\}_{k \leq K_1}) - g_p(1; \{z_k\}_{k \leq K_1})],
 \end{aligned}$$

where  $g_p(\cdot; \{z_k\}_{k \leq K_1})$  is defined in Lemma 5.

By Lemma 5, when  $p = 1$ ,  $g_1(\cdot; \{z_k\}_{k \leq K_1})$  is strictly convex and takes minimum at  $q^* = \frac{1+p}{2} = 1$ , thus

$$g_1(2; \{z_k\}_{k \leq K_1}) - g_1(1; \{z_k\}_{k \leq K_1}) > 0, \quad (75)$$

then we know that

$$\inf_{\mathbf{w}_j \in \mathbb{S}^{D-1}, \cos(\mathbf{w}_j, \bar{\mu}_+) = 1-\delta} \frac{1}{\|\bar{\mu}_+\|} \left\langle \mathcal{P}_{\mathbf{w}_j}^\perp x^{(1)}(\mathbf{w}_j), \bar{\mu}_+ \right\rangle := \Delta_1(\delta) \geq 0. \quad (76)$$

However,  $\Delta_1(\delta)$  can not be zero: If this is the case, since the set  $\{\mathbf{w}_j : \mathbf{w}_j \in \mathbb{S}^{D-1}, \cos(\mathbf{w}_j, \bar{\mu}_+) = 1-\delta\}$  is compact and  $\frac{1}{\|\bar{\mu}_+\|} \left\langle \mathcal{P}_{\mathbf{w}_j}^\perp x^{(1)}(\mathbf{w}_j), \bar{\mu}_+ \right\rangle$  is continuous on this set. It attains minimum 0 at some  $\mathbf{w}_j$ , which implies the non-strong convexity of  $g_1(\cdot; \{z_k\}_{k \leq K_1})$  that, by Lemma 5, requires all  $z_k, k \leq K_1$  to be equal to each other (This can only happen if  $\cos(\mathbf{w}_j, \bar{\mu}_+) = 1$ ). Contradiction. Then one must have

$$\inf_{\mathbf{w}_j \in \mathbb{S}^{D-1}, \cos(\mathbf{w}_j, \bar{\mu}_+) = 1-\delta} \frac{1}{\|\bar{\mu}_+\|} \left\langle \mathcal{P}_{\mathbf{w}_j}^\perp x^{(1)}(\mathbf{w}_j), \bar{\mu}_+ \right\rangle := \Delta_1(\delta) > 0. \quad (77)$$

Similarly, by Lemma 5, when  $p \geq 3$ ,  $g_1(\cdot; \{z_k\}_{k \leq K_1})$  is strictly convex and takes minimum at  $q^* = \frac{1+p}{2} \geq 2$ , thus

$$g_1(2; \{z_k\}_{k \leq K_1}) - g_1(1; \{z_k\}_{k \leq K_1}) < 0, \quad (78)$$

then we know that

$$\inf_{\mathbf{w}_j \in \mathbb{S}^{D-1}, \cos(\mathbf{w}_j, \bar{\boldsymbol{\mu}}_+) = 1-\delta} \frac{1}{\|\bar{\boldsymbol{\mu}}_+\|} \left\langle \mathcal{P}_{\mathbf{w}_j}^\perp x^{(1)}(\mathbf{w}_j), \bar{\boldsymbol{\mu}}_+ \right\rangle := \Delta_p(\delta) \leq 0. \quad (79)$$

Using the same argument, we eliminate the case of this infimum being zero. Then one must have

$$\inf_{\mathbf{w}_j \in \mathbb{S}^{D-1}, \cos(\mathbf{w}_j, \bar{\boldsymbol{\mu}}_+) = 1-\delta} \frac{1}{\|\bar{\boldsymbol{\mu}}_+\|} \left\langle \mathcal{P}_{\mathbf{w}_j}^\perp x^{(1)}(\mathbf{w}_j), \bar{\boldsymbol{\mu}}_+ \right\rangle := \Delta_p(\delta) < 0. \quad (80)$$

The second part concerns about  $\frac{d}{dt} \cos(\mathbf{w}_j(t), \boldsymbol{\mu}_k) \Big|_{\cos(\mathbf{w}_j(t), \boldsymbol{\mu}_k) = 1-\delta}$ , for some  $k \leq K_1$ . Without loss of generality, we let  $k = 1$ . Thus we intend to show  $\frac{d}{dt} \cos(\mathbf{w}_j(t), \boldsymbol{\mu}_1) \Big|_{\cos(\mathbf{w}_j(t), \boldsymbol{\mu}_1) = 1-\delta}$  is positive.

We also let  $\zeta := 1 - (1 - \delta)^2$ , so that the condition  $\cos(\mathbf{w}_j(t), \boldsymbol{\mu}_1) = 1 - \delta$  becomes  $\cos(\mathbf{w}_j(t), \boldsymbol{\mu}_1) = \sqrt{1 - \zeta}$ .

Since  $\sum_{k=1}^K \left| \left\langle \frac{\mathbf{w}_j}{\|\mathbf{w}_j\|}, \boldsymbol{\mu}_k \right\rangle \right|^2 \leq 1$ ,  $\cos(\mathbf{w}_j(t), \boldsymbol{\mu}_1) = \sqrt{1 - \zeta}$  implies  $\sum_{l=2}^K \left| \left\langle \frac{\mathbf{w}_j}{\|\mathbf{w}_j\|}, \boldsymbol{\mu}_l \right\rangle \right|^2 \leq \zeta$ .

Similar to the first part of the proof, we have

$$\frac{d}{dt} \cos(\mathbf{w}_j(t), \boldsymbol{\mu}_1) \quad (81)$$

$$= \left\langle \frac{d}{dt} \frac{\mathbf{w}_j(t)}{\|\mathbf{w}_j(t)\|}, \boldsymbol{\mu}_1 \right\rangle \quad (82)$$

$$= \left\langle \mathcal{P}_{\mathbf{w}_j(t)}^\perp x^{(p)}(\mathbf{w}_j(t)), \boldsymbol{\mu}_1 \right\rangle + \left\langle \frac{d}{dt} \frac{\mathbf{w}_j(t)}{\|\mathbf{w}_j(t)\|} - \text{sign}(v_j(0)) \mathcal{P}_{\mathbf{w}_j(t)}^\perp x^{(p)}(\mathbf{w}_j(t)), \boldsymbol{\mu}_1 \right\rangle, \quad (83)$$

if we can show that for some choice of  $\delta > 0$  (or equivalently some  $\zeta > 0$ ),

$$\inf_{\mathbf{w}_j \in \mathbb{S}^{D-1}, \cos(\mathbf{w}_j, \boldsymbol{\mu}_1) = 1-\delta} \left\langle \mathcal{P}_{\mathbf{w}_j}^\perp x^{(1)}(\mathbf{w}_j), \boldsymbol{\mu}_1 \right\rangle := \Lambda_p(\delta) > 0, \quad (84)$$

then we pick  $\epsilon \leq \epsilon_0 = \frac{\Lambda_p(\delta)}{8\sqrt{h}M^2Kp}$ , and by Lemma 1, we have that for  $\forall t \leq T = \frac{1}{4K} \log \frac{1}{\sqrt{h}\epsilon}$ ,

$$\frac{d}{dt} \cos(\mathbf{w}_j(t), \boldsymbol{\mu}_1) \Big|_{\cos(\mathbf{w}_j(t), \boldsymbol{\mu}_1) = 1-\delta} \geq \Lambda_p(\delta) - \max_j \left\| \frac{d}{dt} \frac{\mathbf{w}_j(t)}{\|\mathbf{w}_j(t)\|} - \text{sign}(v_j(0)) \mathcal{P}_{\mathbf{w}_j(t)}^\perp x^{(p)}(\mathbf{w}_j(t)) \right\| \quad (85)$$

$$\geq \Lambda_p(\delta) - 4\epsilon\sqrt{h}M^2Kp \geq \frac{\Lambda_p(\delta)}{2} > 0, \quad (86)$$

which is what we stated in the theorem. The remaining proof is to find a lower bound on  $\Lambda_p(\delta)$ , which is given by

$$\begin{aligned}
 & \left\langle \mathcal{P}_{\mathbf{w}_j}^{\perp} x^{(p)}(\mathbf{w}_j), \boldsymbol{\mu}_1 \right\rangle \\
 &= \left\langle \left( I - \frac{\mathbf{w}_j \mathbf{w}_j^T}{\|\mathbf{w}_j\|^2} \right) \sum_k^K \gamma_k(\mathbf{w}_j) y_k \mathbf{x}_k \cdot p[\cos(\mathbf{x}_k, \mathbf{w}_j)]^{p-1}, \boldsymbol{\mu}_1 \right\rangle \\
 &= \left\langle \left( I - \frac{\mathbf{w}_j \mathbf{w}_j^T}{\|\mathbf{w}_j\|^2} \right) \sum_{k \leq K_1} \gamma_k(\mathbf{w}_j) y_k \mathbf{x}_k \cdot p[\cos(\mathbf{x}_k, \mathbf{w}_j)]^{p-1}, \boldsymbol{\mu}_1 \right\rangle \\
 &= p \cos^{p-1}(\boldsymbol{\mu}_1, \mathbf{w}_j) (1 - \cos^2(\boldsymbol{\mu}_1, \mathbf{w}_j)) + \sum_{l=2}^K \gamma_l(\mathbf{w}_j) y_l p \cos^p(\boldsymbol{\mu}_l, \mathbf{w}_j) \cos(\mathbf{w}_j, \boldsymbol{\mu}_1) \\
 &= \sqrt{1-\zeta} \left( p(1-\zeta)^{\frac{p-2}{2}} \zeta + \sum_{l=2}^K \gamma_l(\mathbf{w}_j) y_l p \cos^p(\boldsymbol{\mu}_l, \mathbf{w}_j) \right) \\
 &\geq \sqrt{1-\zeta} p \left( (1-\zeta)^{\frac{p-2}{2}} \zeta - \sum_{l=2}^K \left| \left\langle \frac{\mathbf{w}_j}{\|\mathbf{w}_j\|}, \boldsymbol{\mu}_l \right\rangle \right|^p \right) \\
 &\geq \sqrt{1-\zeta} p \left( (1-\zeta)^{\frac{p-2}{2}} \zeta - \left( \sum_{l=2}^K \left| \left\langle \frac{\mathbf{w}_j}{\|\mathbf{w}_j\|}, \boldsymbol{\mu}_l \right\rangle \right|^2 \right)^{\frac{p}{2}} \right) \\
 &\geq \sqrt{1-\zeta} p \left( (1-\zeta)^{\frac{p-2}{2}} \zeta - \zeta^{\frac{p}{2}} \right) \geq \sqrt{1-\zeta} p \left( \left( 1 - \frac{p-2}{2} \zeta \right) \zeta - \zeta^{\frac{p}{2}} \right) = p\zeta + o(\zeta).
 \end{aligned}$$

Therefore, as long as  $\zeta > 0$  is small enough, which can be achieved by picking some  $\delta < \bar{\delta}_2 < 1$ , then  $\inf_{\mathbf{w}_j \in \mathbb{S}^{D-1}, \cos(\mathbf{w}_j, \boldsymbol{\mu}_1) = 1-\delta} \left\langle \mathcal{P}_{\mathbf{w}_j}^{\perp} x^{(1)}(\mathbf{w}_j), \boldsymbol{\mu}_1 \right\rangle = \Lambda_p(\delta)$  is positive.  $\square$

## D.5. Proof for Auxiliary Lemmas

**Balancedness:** Under GF, balancedness (Du et al., 2018) is preserved:  $v_j^2(t) - \|\mathbf{w}_j(t)\|^2 = 0, \forall t \geq 0, \forall j \in [h]$ , from the fact that:

$$\begin{aligned}
 \frac{d}{dt} \|\mathbf{w}_j\|^2 &= \langle \mathbf{w}_j, \dot{\mathbf{w}}_j \rangle \\
 &= -2 \sum_{k=1}^K \gamma_k(\mathbf{w}_j) \nabla_{\hat{y}} \ell(y_k, f_p(\mathbf{x}_k; W, v)) v_j \left( \frac{p[\sigma(\langle \mathbf{x}_k, \mathbf{w}_j \rangle)]^{p-1}}{\|\mathbf{w}_j\|^{p-1}} \langle \mathbf{w}_j, \mathbf{x}_k \rangle - (p-1) \frac{[\sigma(\langle \mathbf{x}_k, \mathbf{w}_j \rangle)]^p}{\|\mathbf{w}_j\|^{p+1}} \|\mathbf{w}_j\|^2 \right) \\
 &= -2 \sum_{k=1}^K \gamma_k(\mathbf{w}_j) \nabla_{\hat{y}} \ell(y_k, f_p(\mathbf{x}_k; W, v)) v_j \left( \frac{p[\sigma(\langle \mathbf{x}_k, \mathbf{w}_j \rangle)]^{p-1}}{\|\mathbf{w}_j\|^p} - (p-1) \frac{[\sigma(\langle \mathbf{x}_k, \mathbf{w}_j \rangle)]^p}{\|\mathbf{w}_j\|^{p-1}} \right) \\
 &= -2 \sum_{k=1}^K \gamma_k(\mathbf{w}_j) \nabla_{\hat{y}} \ell(y_k, f_p(\mathbf{x}_k; W, v)) v_j \frac{[\sigma(\langle \mathbf{x}_k, \mathbf{w}_j \rangle)]^p}{\|\mathbf{w}_j\|^{p-1}} \\
 &= \frac{d}{dt} v_j^2
 \end{aligned}$$

In addition,  $\text{sign}(v_j(t)) = \text{sign}(v_j(0)), \forall t \geq 0, \forall j \in [h]$ , and the dynamical behaviors of neurons will be divided into two types, depending on  $\text{sign}(v_j(0))$ . Therefore, throughout the gradient flow trajectory, we have  $v_j = \text{sign}(v_j(0))\|\mathbf{w}_j\|$ . This fact will be used in the subsequent proof.

*Proof for Lemma 4.* Under gradient flow, we have

$$\frac{d}{dt} \mathbf{w}_j = - \sum_{k=1}^K \gamma_k(\mathbf{w}_j) \nabla_{\hat{y}} \ell(y_k, f_p(\mathbf{x}_k; W, v)) v_j \left( \frac{p[\sigma(\langle \mathbf{x}_k, \mathbf{w}_j \rangle)]^{p-1}}{\|\mathbf{w}_j\|^{p-1}} \mathbf{x}_k - (p-1) \frac{[\sigma(\langle \mathbf{x}_k, \mathbf{w}_j \rangle)]^p}{\|\mathbf{w}_j\|^{p+1}} \mathbf{w}_j \right). \quad (87)$$

and for  $\|\mathbf{w}_j\|$ ,

$$\begin{aligned}
 \frac{d}{dt} \|\mathbf{w}_j\|^2 &= \langle \mathbf{w}_j, \dot{\mathbf{w}}_j \rangle \\
 &= -2 \sum_{k=1}^K \gamma_k(\mathbf{w}_j) \nabla_{\hat{y}} \ell(y_k, f_p(\mathbf{x}_k; W, v)) v_j \left( \frac{p[\sigma(\langle \mathbf{x}_k, \mathbf{w}_j \rangle)]^{p-1}}{\|\mathbf{w}_j\|^{p-1}} \langle \mathbf{w}_j, \mathbf{x}_k \rangle - (p-1) \frac{[\sigma(\langle \mathbf{x}_k, \mathbf{w}_j \rangle)]^p}{\|\mathbf{w}_j\|^{p+1}} \|\mathbf{w}_j\|^2 \right) \\
 &= -2 \sum_{k=1}^K \gamma_k(\mathbf{w}_j) \nabla_{\hat{y}} \ell(y_k, f_p(\mathbf{x}_k; W, v)) v_j \left( \frac{p[\sigma(\langle \mathbf{x}_k, \mathbf{w}_j \rangle)]^{p-1}}{\|\mathbf{w}_j\|^p} - (p-1) \frac{[\sigma(\langle \mathbf{x}_k, \mathbf{w}_j \rangle)]^p}{\|\mathbf{w}_j\|^{p-1}} \right) \\
 &= -2 \sum_{k=1}^K \gamma_k(\mathbf{w}_j) \nabla_{\hat{y}} \ell(y_k, f_p(\mathbf{x}_k; W, v)) v_j \frac{[\sigma(\langle \mathbf{x}_k, \mathbf{w}_j \rangle)]^p}{\|\mathbf{w}_j\|^{p-1}}
 \end{aligned}$$

Balanced initialization enforces  $v_j = \text{sign}(v_j(0))\|\mathbf{w}_j\|$ , hence

$$\frac{d}{dt} \|\mathbf{w}_j\|^2 = -2 \sum_{k=1}^K \gamma_k(\mathbf{w}_j) \nabla_{\hat{y}} \ell(y_k, f_p(\mathbf{x}_k; W, v)) \text{sign}(v_j(0)) \|\mathbf{w}_j\|^2 \frac{[\sigma(\langle \mathbf{x}_k, \mathbf{w}_j \rangle)]^p}{\|\mathbf{w}_j\|^p}. \quad (88)$$

Let  $T := \inf\{t : \max_i |f(\mathbf{x}_k; W(t), v(t))| > 2\epsilon\sqrt{h}M^2\}$ , then  $\forall t \leq T, j \in [h]$ , we have

$$\begin{aligned}
 \frac{d}{dt} \|\mathbf{w}_j\|^2 &= -2 \sum_{k=1}^K \gamma_k(\mathbf{w}_j) \nabla_{\hat{y}} \ell(y_k, f_p(\mathbf{x}_k; W, v)) \text{sign}(v_j(0)) \|\mathbf{w}_j\|^2 \frac{[\sigma(\langle \mathbf{x}_k, \mathbf{w}_j \rangle)]^p}{\|\mathbf{w}_j\|^p} \\
 &= -2 \sum_{k=1}^K \gamma_k(\mathbf{w}_j) \nabla_{\hat{y}} \ell(y_k, f(\mathbf{x}_k; W, v)) \text{sign}(v_j(0)) \|\mathbf{w}_j\|^2 \frac{(\langle \mathbf{x}_k, \mathbf{w}_j \rangle)^p}{\|\mathbf{w}_j\|^p} \\
 &\leq 2 \sum_{k=1}^K |\nabla_{\hat{y}} \ell(y_k, f(\mathbf{x}_k; W, v))| \|\mathbf{w}_j\|^2 \\
 &\leq 2 \sum_{k=1}^K (|y_k| + 2|f(\mathbf{x}_k; W, v)|) \|\mathbf{w}_j\|^2 \\
 &\leq 2 \sum_{k=1}^K (1 + 4\epsilon\sqrt{h}M^2) \|\mathbf{w}_j\| \\
 &\leq 2n(+4\epsilon\sqrt{h}M^2) \|\mathbf{w}_j\|^2.
 \end{aligned} \quad (89)$$

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

$$|f(x_i; W, v)| = \left| \sum_{j \in [h]} \mathbb{1}_{\langle \mathbf{w}_j, \mathbf{x}_k \rangle > 0} v_j \frac{(\langle \mathbf{w}_j, \mathbf{x}_k \rangle)^p}{\|\mathbf{w}_j\|^p} \right| \leq \sum_{j \in [h]} \|\mathbf{w}_j\|^2 \leq h \max_{j \in [h]} \|\mathbf{w}_j\|^2,$$

which implies " $|f(\mathbf{x}_k; W(t), v(t))| > 2\epsilon\sqrt{h}M^2 \Rightarrow \exists j, s.t. \|\mathbf{w}_j(t)\|^2 > \frac{2\epsilon M^2}{\sqrt{h}}$ ".

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

$$\frac{d}{dt} \|w_{j^*}\|^2 \leq 2n(+4\epsilon\sqrt{h}M^2) \|w_{j^*}\|^2. \quad (90)$$

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

$$\begin{aligned}
 \|w_{j^*}(t)\|^2 &\leq \exp\left(2n(+4\epsilon\sqrt{h}M^2)t\right) \|w_{j^*}(0)\|^2, \\
 &= \exp\left(2n(+4\epsilon\sqrt{h}M^2)t\right) \epsilon^2 \|W_0\|_{:,j^*}^2 \\
 &\leq \exp\left(2n(+4\epsilon\sqrt{h}M^2)t\right) \epsilon^2 M^2.
 \end{aligned}$$

Suppose  $\tau_{j^*} < \frac{1}{4n} \log\left(\frac{1}{\sqrt{h}\epsilon}\right)$ , then by the continuity of  $\|w_{j^*}(t)\|^2$ , we have

$$\begin{aligned} \frac{2\epsilon M^2}{\sqrt{h}} &\leq \|w_{j^*}(\tau_{j^*})\|^2 \leq \exp\left(2n(+4\epsilon\sqrt{h}M^2)\tau_{j^*}\right)\epsilon^2 M^2 \\ &\leq \exp\left(2n(+4\epsilon\sqrt{h}M^2)\frac{1}{4n} \log\left(\frac{1}{\sqrt{h}\epsilon}\right)\right)\epsilon^2 M^2 \\ &\leq \exp\left(\frac{1+4\epsilon\sqrt{h}M^2}{2} \log\left(\frac{1}{\sqrt{h}\epsilon}\right)\right)\epsilon^2 M^2 \\ &\leq \exp\left(\log\left(\frac{1}{\sqrt{h}\epsilon}\right)\right)\epsilon^2 M^2 = \frac{\epsilon M^2}{\sqrt{h}}, \end{aligned}$$

which leads to a contradiction  $2\epsilon \leq \epsilon$ . Therefore, one must have  $T \geq \tau_{j^*} \geq \frac{1}{4n} \log\left(\frac{1}{\sqrt{h}\epsilon}\right)$ . This finishes the proof.  $\square$

*Proof of Lemma 5.* Since

$$g'_p(q; \{z_i\}_{i=1}^n) = \left(\sum_{k=1}^K z_i^q \log z_i\right) \left(\sum_{k=1}^K z_i^{p+1-q}\right) - \left(\sum_{k=1}^K z_i^q\right) \left(\sum_{k=1}^K z_i^{p+1-q} \log z_i\right), \quad (91)$$

we immediately find  $g'_p(q^*; \{z_i\}_{i=1}^n) = 0$ . Now we compute the second-order derivative

$$\begin{aligned} g''_p(q; \{z_i\}_{i=1}^n) &= -2 \left(\sum_{k=1}^K z_i^q \log z_i\right) \left(\sum_{k=1}^K z_i^{p+1-q} \log z_i\right) \\ &\quad + \left(\sum_{k=1}^K z_i^q\right) \left(\sum_{k=1}^K z_i^{p+1-q} \log^2 z_i\right) + \left(\sum_{k=1}^K z_i^q \log^2 z_i\right) \left(\sum_{k=1}^K z_i^{p+1-q}\right) \\ &= \sum_{1 \leq i, j \leq n} z_i^q z_j^{p+1-q} (-2 \log z_i \log z_j) \\ &\quad + \sum_{1 \leq i, j \leq n} z_i^q z_j^{p+1-q} \log^2 z_j + \sum_{1 \leq i, j \leq n} z_i^q z_j^{p+1-q} \log^2 z_i \\ &= \sum_{1 \leq i, j \leq n} z_i^q z_j^{p+1-q} (\log z_i - \log z_j)^2 \geq 0, \end{aligned}$$

and the equality holds only when  $z_1 = \dots = z_n$ . The desired results follow.  $\square$