

# THE CONVEX GEOMETRY OF BACKPROPAGATION: NEURAL NETWORK GRADIENT FLOWS CONVERGE TO EXTREME POINTS OF THE DUAL CONVEX PROGRAM

**Yifei Wang**

Department of Electrical Engineering  
Stanford University  
Stanford, CA 94305, USA  
wangyf18@stanford.edu

**Mert Pilanci**

Department of Electrical Engineering  
Stanford University  
Stanford, CA 94305, USA  
pilanci@stanford.edu

## ABSTRACT

We study non-convex subgradient flows for training two-layer ReLU neural networks from a convex geometry and duality perspective. We characterize the implicit bias of unregularized non-convex gradient flow as convex regularization of an equivalent convex model. We then show that the limit points of non-convex subgradient flows can be identified via primal-dual correspondence in this convex optimization problem. Moreover, we derive a sufficient condition on the dual variables which ensures that the stationary points of the non-convex objective are the KKT points of the convex objective, thus proving convergence of non-convex gradient flows to the global optimum. For a class of regular training data distributions such as orthogonal separable data, we show that this sufficient condition holds. Therefore, non-convex gradient flows converge to optimal solutions of a convex optimization problem. We present numerical results verifying the predictions of our theory for non-convex subgradient descent.

## 1 INTRODUCTION

Neural networks (NNs) exhibit remarkable empirical performance in various machine learning tasks. However, a full characterization of the optimization and generalization properties of NNs is far from complete. Non-linear operations inherent to the structure of NNs, over-parameterization and the associated highly nonconvex training problem makes their theoretical analysis quite challenging.

In over-parameterized models such as NNs, one natural question arises: Which particular solution does gradient descent/gradient flow find in unregularized NN training problems? Suppose that  $\mathbf{X} \in \mathbb{R}^{N \times d}$  is the training data matrix and  $\mathbf{y} \in \{1, -1\}^N$  is the label vector. For linear classification problems such as logistic regression, it is known that gradient descent (GD) exhibits implicit regularization properties, see, e.g., (Soudry et al., 2018; Gunasekar et al., 2018). To be precise, under certain assumptions, GD converges to the following solution which maximizes the margin:

$$\arg \min_{\mathbf{w} \in \mathbb{R}^d} \frac{1}{2} \|\mathbf{w}\|_2^2, \text{ s.t. } y_n \mathbf{w}^T \mathbf{x}_n \geq 1, n \in [N]. \quad (1)$$

Here we denote  $[N] = \{1, \dots, N\}$ . Recently, there are several results on the implicit regularization of the (stochastic) gradient descent method for NNs. In (Lyu & Li, 2019), for the multi-layer homogeneous network with exponential or cross-entropy loss, with separable training data, it is shown that the gradient flow (GF) and GD finds a stationary point of the following non-convex max-margin problem:

$$\arg \min_{\boldsymbol{\theta}} \frac{1}{2} \|\boldsymbol{\theta}\|_2^2, \text{ s.t. } y_n f(\boldsymbol{\theta}; \mathbf{x}_n) \geq 1, n \in [N], \quad (2)$$

where  $f(\boldsymbol{\theta}; \mathbf{x})$  represents the output of the neural network with parameter  $\boldsymbol{\theta}$  given input  $\mathbf{x}$ . In (Phuong & Lampert, 2021), by further assuming the orthogonal separability of the training data, it is shown that all neurons converge to one of the two max-margin classifiers. One corresponds to the data with positive labels, while the other corresponds to the data with negative labels.

However, as the max-margin problem of the neural network (2) is a non-convex optimization problem, the existing results only guarantee that it is a stationary point which can be a local minimizer or even a saddle point. In other words, the global optimality is not guaranteed.

In a different line of work (Pilanci & Ergen, 2020; Ergen & Pilanci, 2020; 2021b), exact convex optimization formulations of two and three-layer ReLU NNs are developed, which have global optimality guarantees in polynomial-time when the data has a polynomial number of hyperplane arrangements, e.g., in any fixed dimension or with convolutional networks of fixed filter size. The convex optimization framework was extended to vector output networks (Sahiner et al., 2021b), quantized networks (Bartan & Pilanci, 2021b), autoencoders (Sahiner et al., 2021c; Gupta et al., 2021), networks with polynomial activation functions (Bartan & Pilanci, 2021a), networks with batch normalization (Ergen et al., 2021), univariate deep ReLU networks, deep linear networks (Ergen & Pilanci, 2021c) and Generative Adversarial Networks (Sahiner et al., 2021a).

In this work, we first derive an equivalent convex program corresponding to the maximal margin problem (2). We then consider non-convex subgradient flow for unregularized logistic loss. We show that the limit points of non-convex subgradient flow can be identified via primal-dual correspondence in the convex optimization problem. We then present a sufficient condition on the dual variable to ensure that all stationary points of the non-convex max-margin problem are KKT points of the convex max-margin problem. For certain regular datasets including orthogonal separable data, we show that this sufficient condition on the dual variable holds, thus implies the convergence of gradient flow on the unregularized problem to the global optimum of the non-convex maximal margin problem (2). Consequently, this enables us to fully characterize the implicit regularization of unregularized gradient flow or gradient descent as convex regularization applied to a convex model.

## 1.1 RELATED WORK

There are several works studying the property of two-layer ReLU networks trained by gradient descent/gradient flow dynamics. The following papers study the gradient descent like dynamics in training two-layer ReLU networks for regression problems. Ma et al. (2020) show that for two-layer ReLU networks, only a group of a few activated neurons dominate the dynamics of gradient descent. In (Mei et al., 2018), the limiting dynamics of stochastic gradient descent (SGD) is captured by the distributional dynamics from a mean-field perspective and they utilize this to prove a general convergence result for noisy SGD. Li et al. (2020) focus on the case where the weights of the second layer are non-negative and they show that the over-parameterized neural network can learn the ground-truth network in polynomial time with polynomial samples. In (Zhou et al., 2021), it is shown that mildly over-parameterized student network can learn the teacher network and all student neurons converge to one of the teacher neurons.

Beyond (Lyu & Li, 2019) and (Phuong & Lampert, 2021), the following papers study the classification problems. In (Chizat & Bach, 2018), under certain assumptions on the training problem, with over-parameterized model, the gradient flow can converge to the global optimum of the training problem. For linear separable data, utilizing the hinge loss for classification, Wang et al. (2019) introduce a perturbed stochastic gradient method and show that it can attain the global optimum of the training problem. Similarly, for linear separable data, Yang et al. (2021) introduce a modified loss based on the hinge loss to enable (stochastic) gradient descent find the global minimum of the training problem, which is also globally optimal for the training problem with the hinge loss.

## 1.2 PROBLEM SETTING

We focus on two-layer neural networks with ReLU activation, i.e.,  $f(\cdot; X) = (XW_1)_+ w_2$ , where  $W_1 \in \mathbb{R}^{d \times m}$ ,  $w_2 \in \mathbb{R}^m$  and  $\theta = (W_1; w_2)$  represents the parameter. Due to the ReLU activation, this neural network is homogeneous, i.e., for any scalar  $c > 0$ , we have  $f(c\theta; X) = c^2 f(\theta; X)$ . The training problem is given by

$$\min_{\theta} \sum_{n=1}^N \ell(y_n f(\theta; x_n)); \quad (3)$$

where  $\ell(q) : \mathbb{R} \rightarrow \mathbb{R}_+$  is the loss function. We focus on the logistic, i.e., cross-entropy loss, i.e.,  $\ell(q) = \log(1 + \exp(-q))$ .

Then, we briefly review gradient descent and gradient flow. The gradient descent takes the update rule

$$\theta(t+1) = \theta(t) - \eta g(t);$$

where  $g(t) \in \partial L(\theta(t))$  and  $\partial$  represents the Clarke's subdifferential.

The gradient flow can be viewed as the gradient descent with infinitesimal step size. The trajectory of the parameter during training is an arc  $\gamma : [0, +\infty) \rightarrow \mathbb{R}^d$ , where  $\dot{\gamma} = -f(\gamma) = -(W_1 w_2)$ ,  $W_1 \in \mathbb{R}^{d \times m}$ ;  $W_2 \in \mathbb{R}^{m \times g}$ . More precisely, the gradient flow is given by the differential inclusion

$$\frac{d}{dt} \theta(t) \in \partial L(\theta(t)); \text{ for } t \geq 0, \text{ almost everywhere.}$$

## 2 MAIN RESULTS

In this section, we present our main results and defer the detailed analysis to the following sections.

Consider the more general multi-class version of the problem with  $K$  classes. Suppose that  $\mathbf{y} \in [K]^N$  is the label vector. Let  $\mathbf{Y} = (y_{n,k})_{n \in [N], k \in [K]} \in \mathbb{R}^{N \times K}$  be the encoded label matrix such that

$$y_{n,k} = \begin{cases} 1; & \text{if } y_n = k; \\ 0; & \text{otherwise} \end{cases}$$

Similarly, we consider the following two-layer vector-output neural networks with ReLU activation:

$$F(\mathbf{x}; X) = \begin{bmatrix} f_1(\mathbf{x}; X) \\ \vdots \\ f_K(\mathbf{x}; X) \end{bmatrix} = \begin{bmatrix} (XW_1^{(1)})_+ w_2^{(1)} \\ \vdots \\ (XW_1^{(K)})_+ w_2^{(K)} \end{bmatrix};$$

where we write  $\mathbf{w}_2 = (w_2^{(1)}; \dots; w_2^{(K)})$ . For  $k = 1; \dots; K$ , we have  $w_2^{(k)} = (W_1^{(k)}; w_2^{(k)})$  where  $W_1^{(k)} \in \mathbb{R}^{N \times m}$  and  $w_2^{(k)} \in \mathbb{R}^m$ . One can view each of the outputs of  $F(\mathbf{x}; X)$  as the output of a two-layer scalar-output neural network. Consider the following training problem:

$$\min_{X \in \mathbb{R}^{N \times m}} \sum_{k=1}^K \sum_{n=1}^N (y_{n,k} - f_k(\mathbf{x}_n; X))^2; \quad (4)$$

According to (Lyu & Li, 2019), the gradient flow and the gradient descent finds a stationary point of the following non-convex max-margin problem:

$$\arg \min_{X \in \mathbb{R}^{N \times m}} \sum_{k=1}^K \frac{1}{2} \|k\|_2^2; \text{ s.t. } y_{n,k} - f_k(\mathbf{x}_n; X) \geq 1; n \in [N]; k \in [K]; \quad (5)$$

Denote the set of all possible hyperplane arrangement as

$$\mathcal{P} = \{ \text{diag}(I(Xw - 0)) \mid w \in \mathbb{R}^d \}; \quad (6)$$

and let  $p = |P|$ . We can also write  $\mathcal{P} = \{D_1; \dots; D_p\}$ . From (Cover, 1965), we have an upper bound  $p \leq 2^r \frac{e(N-1)}{r}$  where  $r = \text{rank}(X)$ . We first reformulate (5) as convex optimization.

**Proposition 1** The non-convex problem (5) is equivalent to the following convex program

$$\begin{aligned} \min_{X \in \mathbb{R}^{N \times m}} & \sum_{k=1}^K \sum_{j=1}^p (k u_{j,k} \|k\|_2 + k u_{j,k}^0 \|k\|_2); \\ \text{s.t. } & \text{diag}(y_k) = \sum_{j=1}^p D_j X (u_{j,k} - u_{j,k}^0) \geq \mathbf{1}; \\ & (2D_j - I)X u_{j,k} \leq 0; (2D_j - I)X u_{j,k}^0 \geq 0; j \in [p]; k \in [K]; \end{aligned} \quad (7)$$

where  $y_k$  is the  $k$ -th column of  $\mathbf{Y}$ . The dual problem of (7) is given by

$$\begin{aligned} \max & \text{tr}(Y^T \mathbf{1}); \\ \text{s.t. } & \text{diag}(y_k) \leq k; \max_{k \in [K]} \sum_{j=1}^p \frac{1}{k} (X^T w)_j \geq 1; k \in [K]; \end{aligned} \quad (8)$$

where  $y_k$  is the  $k$ -th column of  $\mathbf{Y}$ .

We present the detailed derivation of the convex formulation (7) and its dual problem (8) in the appendix. Given  $u \in \mathbb{R}^d$ , we define  $D(u) = \text{diag}(1(Xu > 0))$ . For two vectors  $u, v \in \mathbb{R}^d$ , we define the cosine angle between  $u$  and  $v$  by  $\cos(u; v) = \frac{u^T v}{\|u\|_2 \|v\|_2}$ .

## 2.1 OUR CONTRIBUTIONS

The following theorem illustrates that for neurons satisfying  $\text{sign}(y_k^T (Xw_{1;i}^{(k)}))_+ = \text{sign}(w_{2;i}^{(k)})$  at initialization,  $w_{1;i}^{(k)}$  align to the direction of  $X^T D(w_{1;i}^{(k)}) y_k$  at a certain time  $T$ , depending on  $\text{sign}(w_{2;i}^{(k)})$  at initialization. In Section 2.3, we show that these are dual extreme points of (7).

**Theorem 1** Consider the  $K$ -class classification training problem (4) for any dataset. Suppose that the neural network is scaled at initialization such that  $\|w_{1;i}^{(k)}\|_2 = \|w_{2;i}^{(k)}\|_2$  for  $i \in [m]$  and  $k \in [K]$ . Assume that at initialization, for  $k \in [K]$ , there exists neuron  $(w_{1;i_k}^{(k)}; w_{2;i_k}^{(k)})$  such that

$$\text{sign}(y_k^T (Xw_{1;i_k}^{(k)}))_+ = \text{sign}(w_{2;i_k}^{(k)}) = s; \quad (9)$$

where  $s \in \{-1, 1\}$ . Consider the subgradient flow applied to the non-convex problem (4). Let  $\epsilon \in (0, 1)$ . Suppose that the initialization is sufficiently close to the origin. Then, for  $k \in [K]$ , there exist  $T = T(\epsilon; k)$  such that

$$\cos(w_{1;i_k}^{(k)}(T); sX^T D(w_{1;i_k}^{(k)}(T)) y_k) \geq 1 - \epsilon.$$

Next, we impose conditions on the dataset to prove a stronger global convergence results on the flow. We say that the dataset  $(X; y)$  is orthogonal separable among multiple classes if for all  $n \in [N]$ ,

$$\begin{aligned} x_n^T x_{n^0} &> 0; \text{ if } y_n = y_{n^0}; \\ x_n^T x_{n^0} &\leq 0; \text{ if } y_n \neq y_{n^0}. \end{aligned}$$

For orthogonal separable dataset among multiple classes, the subgradient flow for the non-convex problem (4) can find the global optimum of (5) up to a scaling constant.

**Theorem 2** Suppose that  $(X; y) \in \mathbb{R}^{N \times d} \times [K]^N$  is orthogonal separable among multiple classes. Consider the non-convex subgradient flow applied to the non-convex problem (4). Suppose that the initialization is sufficiently close to the origin and scaled as in Theorem 1. Then, the non-convex subgradient flow converges to the global optimum of the convex problem (7) and hence the non-convex objective (5) up to scaling.

Therefore, the above result characterizes the implicit regularization of unregularized gradient flow as convex regularization, i.e., group  $\ell_1$  norm, in the convex formulation (7). It is remarkable that group sparsity is enforced by small initialization magnitude with no explicit form of regularization.

## 2.2 CONVEX GEOMETRY OF NEURAL GRADIENT FLOW

Suppose that  $X \in \mathbb{R}^{N \times d}$ . Here we provide an interesting geometric interpretation behind the formula

$$\cos(u; X^T D(u)) \geq 1 - \epsilon,$$

which describes dual extreme points to which hidden neurons approach to as predicted by Theorem 1. We now explain the geometric intuition behind this result. Consider an ellipsoid  $\mathcal{E} = \{Xu : \|u\|_2 \leq 1\}$ . A positive extreme point of this ellipsoid along the direction is defined by  $\arg \max_{u: \|u\|_2 \leq 1} u^T Xu$ , which is given by the formula  $\frac{X^T X u}{\|X^T X u\|_2}$ . Next, we consider the rectified ellipsoid set  $\mathcal{Q} := \{f(Xu)_+ : \|u\|_2 \leq 1\}$  introduced in (Ergen & Pilanci, 2021a) and shown in Figure 1. The constraint  $\max_{u: \|u\|_2 \leq 1} j^T (Xu)_+ \geq 1$  on  $\mathcal{Q}$  is equivalent to  $2 \in \mathcal{Q}$ . Here  $\mathcal{Q}$  is the absolute polar set  $\mathcal{Q}^*$ , which appears as a constraint in the convex problem (8) and is defined as the following convex set

$$\mathcal{Q} = \{f : \max_{z \in \mathcal{Q}^*} j^T z \geq 1\}; \quad (10)$$

An extreme point of this non-convex body along the directions given by the solution of the problem

$$\max_{u: \|u\|_2 = 1} \sum_{j=1}^P (Xu)_j = \max_{D_j \in \mathcal{P}} \max_{u: \|u\|_2 = 1; (2D_j - I)Xu \geq 0} \sum_{j=1}^P D_j Xu : \quad (11)$$

Here,  $(\lambda; u)$  are primal-dual pairs as they appear in the convex dual problem. First, note that a stationary point of gradient flow on the objective (11) is given by the identity  $\sum_{j=1}^P (Xu)_j = c$ , where  $c$  is a constant. In particular, by picking the zero as the subgradient  $(\lambda_n^T u)_+$  when  $x_n^T u = 0$ ,

$$u = \frac{X^T D(u)}{\|X^T D(u)\|_2} = \frac{\sum_{n=1}^P \sum_{j=1}^N x_n I(u^T x_n > 0)}{\|\sum_{n=1}^P \sum_{j=1}^N x_n I(u^T x_n > 0)\|_2} : \quad (12)$$

Note that the formula  $\cos(\lambda; X^T D(u)) > 1$  appearing in Theorem 1 shows that gradient flow reaches the extreme points of projected ellipsoids  $\{Xu : \|u\|_2 = 1\}$  in the direction of  $\lambda = y_k$ , where  $D_j \in \mathcal{P}$  corresponds to a valid hyperplane arrangement. This interesting phenomenon is depicted in Figures 3 and 4. The one-dimensional spikes in Figures 1 and 3 are projected ellipsoids. Detailed setup for Figure 1 to 4 and additional experiments can be found in Appendix F.

Figure 1: Rectified Ellipsoid  $Q := \sum_{j=1}^P (Xu)_j : \|u\|_2 = 1$  and its extreme points (spikes).

Figure 2: Convex absolute polar set of the Rectified Ellipsoid (purple) and other dual constraints (grey).

Figure 3: Trajectories of  $(Xw_{1;j})_+$  along the training dynamics of gradient descent.

Figure 4: Trajectories of  $w_{1;j} = \frac{w_{1;j}}{\|w_{1;j}\|_2}$  along the training dynamics of gradient descent.

Figure 5: Two-layer ReLU network gradient descent dynamics on an orthogonal separable dataset.  $w_{1;j} = \frac{w_{1;j}}{\|w_{1;j}\|_2}$  is the normalized vector of the  $j$ th hidden neuron in the first layer.

### 3 CONVEX MAX-MARGIN PROBLEM

In this section, we consider the equivalent convex model of the max-margin problem and its optimality conditions. We primarily focus on the binary classification problem for simplicity, which are later extended to the multi-class case. We can reformulate the nonconvex max-margin problem (2) as

$$\min \frac{1}{2} (k w_1 k_F^2 + k w_2 k_2^2); \text{ s.t. } Y (X w_1)_+ + w_2 \leq 1; \quad (13)$$

where  $Y = \text{diag}(y)$ . This is a nonconvex optimization problem due to the ReLU activation and the two-layer structure of neural network. Analogous to the convex formulation introduced in (Pilanci & Ergen, 2020) for regularized training problem of neural network, we can provide a convex optimization formulation of (13) and derive the dual problem.

Proposition 2 The problem (13) is equivalent to

$$\begin{aligned} P_{\text{cvx}} = \min_{\{u_j\}_{j=1}^m} & \sum_{j=1}^m (k u_j k_2 + k u_j^0 k_2); \\ \text{s.t. } & Y \sum_{j=1}^m D_j X (u_j^0 - u_j) \leq 1; \\ & (2D_j - I) X u_j \leq 0; (2D_j - I) X u_j^0 \leq 0; \forall j \in [m]; \end{aligned} \quad (14)$$

The dual problem of (14) is given by

$$D = \max y^T \text{ s.t. } Y \leq 0; \max_{u: k u k_2 \leq 1} \sum_{j=1}^m j^T (X^T u)_+ \leq 1; \quad (15)$$

The following proposition gives a characterization of the KKT point of the non-convex max-margin problem (2). The definition of B-subdifferential can be found in Appendix A.

Proposition 3 Let  $(w_1; w_2; \lambda)$  be a KKT point of the non-convex max-margin problem (2) (in terms of B-subdifferential). Suppose that  $\lambda_{i,j} \in \{0, 1\}$  for certain  $i \in [m]$ . Then, there exists a diagonal matrix  $\hat{D}_i \in \mathbb{R}^{n \times n}$  satisfying

$$\begin{aligned} (\hat{D}_i)_{n,n} &= 1; \text{ for } x_n^T w_{1,i} > 0; \\ (\hat{D}_i)_{n,n} &\in [0, 1]; \text{ for } x_n^T w_{1,i} = 0; \\ (\hat{D}_i)_{n,n} &= 0; \text{ for } x_n^T w_{1,i} < 0; \end{aligned}$$

such that

$$\frac{w_{1,i}}{w_{2,i}} = X^T \hat{D}_i; \quad k X^T \hat{D}_i k_2 = 1;$$

Based on the characterization of the KKT point of the non-convex max-margin problem (2), we provide an equivalent condition to ensure that it is also the KKT point of the convex max-margin problem (14).

Theorem 3 The KKT point of the non-convex max-margin problem (2) (in terms of B-subdifferential) corresponds to a KKT point of the convex max-margin problem (14) if  $\lambda$  is dual feasible, i.e.,

$$\max_{u: k u k_2 \leq 1} \sum_{j=1}^m j^T (X u)_+ \leq 1; \quad (16)$$

This condition is equivalent to for all  $D_j \in \mathcal{P}$ , the dual variable  $\lambda$  satisfies that

$$\max_{k u k_2 \leq 1; (2D_j - I) X u \leq 0} \sum_{j=1}^m j^T D_j X u \leq 1; \quad (17)$$

#### 3.1 DUAL FEASIBILITY OF THE DUAL VARIABLE

A natural question arises: is it possible to examine whether  $\lambda$  is feasible in the dual problem? We say the dataset  $(X; y)$  is orthogonal separable if for all  $n^0 \in [N]$ ,

$$\begin{aligned} x_n^T x_{n^0} &> 0; \text{ if } y_n = y_{n^0}; \\ x_n^T x_{n^0} &\leq 0; \text{ if } y_n \neq y_{n^0}; \end{aligned}$$

For orthogonal separable data, as long as the induced diagonal matrices in Proposition 3 cover the positive part and the negative part of the labels, the KKT point of the non-convex max-margin problem (2) is the KKT point of the convex max-margin problem (14).

Proposition 4 Suppose that  $(X; y)$  is orthogonal separable. Suppose that the KKT point of the non-convex problem include two neurons  $(w_{1;i_+}; w_{2;i_+})$  and  $(w_{1;i_-}; w_{2;i_-})$  such that the corresponding diagonal matrices  $\hat{D}_{i_+}$  and  $\hat{D}_{i_-}$  defined in Proposition 3 satisfy that

$$\hat{D}_{i_+} = \text{diag}(I(y = 1)); \quad \hat{D}_{i_-} = \text{diag}(I(y = -1));$$

Then, the dual variable is dual feasible, i.e., satisfying (6).

The spike-free matrices discussed in (Ergen & Pilanci, 2021a) also makes examining the dual feasibility easier. The definition of spike-free matrices can be found in Appendix A

Proposition 5 Suppose that  $X$  is spike-free. Suppose that the KKT point of the non-convex problem include two neurons  $(w_{1;i_+}; w_{2;i_+})$  and  $(w_{1;i_-}; w_{2;i_-})$  such that the corresponding diagonal matrices  $\hat{D}_{i_+}$  and  $\hat{D}_{i_-}$  defined in Proposition 3 satisfy that

$$\hat{D}_{i_+} = \text{diag}(I(y = 1)); \quad \hat{D}_{i_-} = \text{diag}(I(y = -1));$$

Then, the dual variable is dual feasible, i.e., satisfying (6).

Remark 1 For the spike-free data, the constraint on the dual problem is equivalent to

$$\begin{aligned} \max_{X_u \geq 0; \|k\|_2 = 1} \sum_j X_{uj}^T X_{uj} - 1; \quad \text{or equivalently} \\ \max_{X_u \geq 0; \|k\|_2 = 1} \sum_j Y_{+j} X_{uj} - 1; \quad \min_{X_u \geq 0} \sum_j Y_{-j} X_{uj} - 1; \end{aligned}$$

#### 4 SUB-GRADIENT FLOW DYNAMICS OF LOGISTIC LOSS

In this section, we consider the following sub-gradient flow of the logistic loss (3)

$$\begin{aligned} \frac{\partial}{\partial t} w_{1;i}(t) &= w_{2;i}(t) \sum_{n: (w_{1;i}(t))^T x_n > 0} \tilde{y}_n(t) x_n(t); \\ \frac{\partial}{\partial t} w_{2;i}(t) &= \sum_{n=1}^N \tilde{y}_n(t) ((w_{1;i}(t))^T x_n(t))_+; \end{aligned} \tag{18}$$

where the  $n$ -th entry of  $\tilde{y}(t) \in \mathbb{R}^N$  is defined

$$\tilde{y}_n = y_n \cdot \sigma(q_n); \quad q_n = y_n (x_n^T W_{1+} + w_{2-}); \tag{19}$$

For simplicity, we omit the term  $\tilde{y}(t)$ . For instance, we write  $w_{1;i} = w_{1;i}(t)$ . To be specific, when  $w_{1;i}^T x_n = 0$ , we select  $\tilde{y}_n$  as the subgradient of  $\sigma(w_{1;i}^T x_n)_+$  with respect to  $w_{1;i}$ . Denote  $u_i = \text{sign}(X u_i)$ . For  $\tilde{y} \in \{-1, 1\}^N$ , we define

$$g(\tilde{y}; e) = \sum_{n: u_n > 0} \tilde{y}_n x_n; \tag{20}$$

For simplicity, we also write

$$g(u; e) := g(\text{sign}(X u); \tilde{y}) = \sum_{n: w_{1;i}^T x_n > 0} \tilde{y}_n x_n; \tag{21}$$

Then, we can rewrite sub-gradient flow of the logistic loss (3) as follows:

$$\frac{\partial}{\partial t} w_{i;1} = w_{2;i} g(u; e); \quad \frac{\partial}{\partial t} w_{i;2} = w_{1;i}^T g(u; e); \tag{22}$$

Assume that the neural network is scaled at initialization,  $\|w_{1;i}(0)\|_2^2 = w_{2;i}^2(0)$  for  $i \in [m]$ . Then, the neural network is scaled for 0.

Lemma 1 Suppose that  $\|w_{1;i}(0)\|_2 = \|w_{2;i}(0)\|_2 > 0$  for  $i \in [m]$ . Then, for any  $\epsilon > 0$ , we have  $\|w_{1;i}(t)\|_2 = \|w_{2;i}(t)\|_2 > 0$ .

According to Lemma 1, for all  $t \geq 0$ ,  $\text{sign}(w_{2;i}(t)) = \text{sign}(w_{2;i}(0))$ . Therefore, we can simply write  $s_i = s_i(t) = \text{sign}(w_{2;i}(t))$ . As the neural network is scaled for  $\epsilon > 0$ , it is interesting to study the dynamics of  $w_{1;i}$  in the polar coordinate. We write  $w_{1;i}(t) = e^{r_i(t)} u_i(t)$ , where  $\|u_i(t)\|_2 = 1$ . The gradient flow in terms of polar coordinate writes

$$\frac{\partial}{\partial t} r_i = s_i u_i^T g(u_i; e); \quad \frac{\partial}{\partial t} u_i = s_i \nabla g(u_i; e) - u_i^T \nabla g(u_i; e) u_i \quad (23)$$

Let  $x_{\max} = \max_{i \in [n]} \|x_i\|_2$ . Define  $g_{\min}$  to be

$$g_{\min} = \min_{\|y\|_2 \leq x_{\max}} \langle \nabla g(y; e), y \rangle; \quad \text{s.t. } \langle \nabla g(y; e), y \rangle \leq 0; \quad \text{where we denote} \quad (24)$$

$$Q = \{y \in \mathbb{R}^d; \|y\|_2 \leq x_{\max}, \langle \nabla g(y; e), y \rangle \leq 0\} \quad (25)$$

As the set  $Q$  is finite, we note that  $g_{\min} > 0$ . We note that when  $\max_{i \in [n]} \|x_i\|_2 \leq \frac{g_{\min}}{4}$ , we have  $\frac{d}{dt} \|u_i(t)\|_2 \leq \frac{g_{\min}}{4}$ . The following lemma shows that with initializations sufficiently close to  $Q$ ,  $\|g(u(t); e(t)) - \nabla g(u(t); e(t))\|_2$  and  $\frac{d}{dt} \|u_i(t)\|_2$  can be very small.

Lemma 2 Suppose that  $\epsilon > 0$  and  $\delta > 0$ . Suppose that  $(u(t); r(t))$  follows the gradient flow (23) with  $s = 1$  and the initialization  $u(0) = u_0$  and  $r(0) = r_0$ . Suppose that  $t_0$  is sufficiently small. Then, the following two statements hold.

- For all  $t \in [t_0, T]$ , we have  $\|g(u(t); e(t)) - \nabla g(u(t); e(t))\|_2 \leq \frac{g_{\min}}{8}$ ;
- For  $t \in [t_0, T]$  such that  $\text{sign}(Xu(t))$  is constant in a small neighborhood of  $t$ , we have  $\frac{d}{dt} \|u_i(t)\|_2 \leq \frac{g_{\min}^2}{16}$ ;

Based on the above lemma on the property of  $(u(t); e(t))$ , we introduce the following lemma to upper-bound the time such that  $\text{sign}(Xu(t))$  approaches or  $\text{sign}(Xu(t))$  changes.

Lemma 3 Let  $c \in (0, 1)$ . Suppose that  $t_0$  satisfies that  $\|u_0\|_2 = 1$  and  $e(0)^T (Xu_0)_+ > 0$ . Suppose that  $(u(t); r(t))$  follows the gradient flow (23) with  $s = 1$  and the initialization  $u(0) = u_0$  and  $r(0) = r_0$ . Let  $v(t) = \frac{g(u(t); e(t))}{\|g(u(t); e(t))\|_2}$ . We write  $v_0 = v(0)$ ,  $v_0 = (v_0)_+$  and  $g_0 = \|g(u_0; e_0)\|_2$ . Denote

$$T = \frac{1}{2g_0} \log \frac{1 + \sqrt{1 + 8 + 1}}{1 + \sqrt{1 + 8 + v_0^T u_0}} \quad (26)$$

For  $c \in (0, 1)$ , define

$$T^{\text{shift}}(c) = \frac{1}{2g_0} \log \frac{1 + \sqrt{1 + 8 + c}}{1 + \sqrt{1 + 8 + v_0^T u_0}} \quad (27)$$

Suppose that  $t_0$  is sufficiently small such that the statements in Lemma 2 hold for  $T$ . Then, at least one of the following event happens

- There exists a time  $\bar{t} \in [0, T]$  such that we have  $\text{sign}(Xu(\bar{t})) = \text{sign}(Xu_0)$  for  $t \in [0, \bar{t}]$  and  $\text{sign}(Xu(\bar{t})) \neq \text{sign}(Xu_0)$ . Let  $u_1 = u(\bar{t})$  and  $v_1 = \lim_{t \downarrow \bar{t}} v(t)$ . If  $u_1^T v_1 > 1$ , then the time  $\bar{t}$  satisfies that  $\bar{t} \leq T^{\text{shift}}(v_1^T u_1)$ . Otherwise, there exists a time  $\bar{t}^0$  satisfying  $\bar{t}^0 \leq T$ ; such that we have  $\text{sign}(Xu(\bar{t}^0)) = \text{sign}(Xu_0)$  for  $t \in [0, \bar{t}^0]$  and  $u(\bar{t}^0)^T v(\bar{t}^0) > 1$ .
- There exists a time  $\bar{t} \in [0, T]$ ; such that we have  $\text{sign}(Xu(\bar{t})) = \text{sign}(Xu_0)$  for  $t \in [0, \bar{t}]$  and  $u(\bar{t})^T v(\bar{t}) > 1$ .



Corollary 1 Suppose that there exists a time  $\bar{T}$  such that we have  $\text{sign}(Xu(t)) = \text{sign}(Xu_0)$  for  $t \in [0; T]$  and  $\text{sign}(Xu(t)) \neq \text{sign}(Xu_0)$ . If  $T > T^{\text{shift}}(v_1^T u_1) = \frac{p-1}{g_0} \log \frac{1+8+v_1^T u_1}{1-8+v_1^T u_1} \log \frac{1+8+v_0^T u_0}{1-8+v_0^T u_0}$ , then, we have  $v_1^T v_1 > 1$ .

Proposition 6 Consider the sub-gradient  $\alpha$  (23) with  $s = 1$  and the initialization  $u(0) = u_0$  and  $r(0) = r_0$ . Here at initialization the neuron  $w_0$  satisfies that  $\|w_0\|_2 = 1$  and  $y^T(Xu_0)_+ > 0$ . Let  $v(t) = \frac{g(u(t); e(t))}{kg(u(t); e(t))_2}$ . For any  $\epsilon > 0$ , for sufficiently small  $r_0$ , there exists a time  $\bar{T} = O(\log(\frac{1}{\epsilon}))$  such that  $u(T)^T v(T) = 1$  and  $\cos(u(T); g(u(T); y)) = 1$ .

Remark 2 The statement of proposition is similar to Lemma 4 in (Maennel et al., 2018). However, their proof contains a problem because they did not consider the change of  $\text{sign}(Xw)$  along the gradient flow. Our proof in Appendix D.4 corrects this error.

We next study the properties of orthogonal separable datasets. Definition 2  $R^d : \|w\|_2 = 1$ g. The following lemma give a sufficient condition on  $w$  to satisfy the condition in Proposition 4.

Lemma 4 Assume that  $(X; y)$  is orthogonal separable. Suppose that  $w \in B$  is a local maximizer of  $y^T(Xw)_+$  in  $B$  and  $(Xw)_+ \neq 0$ . Then,  $\langle w; x_n \rangle > 0$  for  $n \in [N]$  such that  $y_n = 1$ . Suppose that  $w \in B$  is a local minimizer of  $y^T(Xw)_+$  in  $B$  and  $(Xw)_+ \neq 0$ . Then,  $\langle w; x_n \rangle > 0$  for  $n \in [N]$  such that  $y_n = -1$ .

We show an equivalent condition of  $B$  being the local maximizer/minimizer of  $y^T(Xu)_+$  in  $B$ .

Proposition 7 Assume that  $(X; y)$  is orthogonal separable. Then,  $B$  is a local maximizer of  $y^T(Xu)_+$  in  $B$  is equivalent to  $\cos(u; g(u; y)) = 1$ . Similarly,  $u \in B$  is a local minimizer of  $y^T(Xu)_+$  in  $B$  is equivalent to  $\cos(u; g(u; y)) = -1$ .

Based on Proposition 4 and 7, we present the main theorem.

Theorem 4 Suppose that the dataset is orthogonal separable and follows the gradient flow. Suppose that the neural network is scaled at initialization,  $\|w_{1;i}(0)\|_2 = |w_{2;i}(0)|$  for all  $i \in [m]$ . For almost all initializations which are sufficiently close to zero, the limiting point of  $\frac{f(t)}{\|k(t)\|_2}$  is  $\frac{f}{k-k_2}$ , where  $f$  is a global minimizer of the max-margin problem (2).

We present a sketch of the proof. According to Proposition 6, for initialization sufficiently close to zero, there exist two neurons and time  $\bar{T} > 0$  such that  $\cos(w_{1;i_+}(\bar{T}); g(w_{1;i_+}(\bar{T}); y)) = 1$  and  $\cos(w_{1;i_-}(\bar{T}); g(w_{1;i_-}(\bar{T}); y)) = -1$ . This implies that  $w_{1;i_+}(\bar{T})$  and  $w_{1;i_-}(\bar{T})$  are sufficiently close to certain stationary points of gradient flow maximizing/minimizing  $y^T(Xu)_+$  over  $B$ , i.e.,  $f \in B$   $\cos(u; g(u; y)) = 1$ g. As the dataset is orthogonal separable, from Proposition 7 and Lemma 4, the induced masking matrices  $D_{i_+}(\bar{T})$  and  $D_{i_-}(\bar{T})$  by  $w_{1;i_+}(\bar{T})/w_{1;i_-}(\bar{T})$  in Proposition 3 satisfy that  $D_{i_+}(\bar{T}) = \text{diag}(I(y = 1))$  and  $D_{i_-}(\bar{T}) = \text{diag}(I(y = -1))$ . According to Lemma 3 in (Phuong & Lampert, 2021), for  $\max\{\bar{T}; T\} \geq g$ , we also have  $D_{i_+}(t) = \text{diag}(I(y = 1))$  and  $D_{i_-}(t) = \text{diag}(I(y = -1))$ . According to Theorem 3 and Proposition 4, the KKT point of the non-convex problem (2) that gradient flow converges to corresponds to the KKT point of the convex problem (14).

## 5 CONCLUSION

We provide a convex formulation of the non-convex max-margin problem for two-layer ReLU neural networks and uncover a primal-dual extreme point relation between non-convex subgradient flow. Under the assumptions on the training data, we show that flows converge to KKT points of the convex max-margin problem, hence a global optimum of the non-convex objective.

## 6 ACKNOWLEDGEMENTS

This work was partially supported by the National Science Foundation under grants ECCS-2037304, DMS-2134248, and US Army Research Office.

## REFERENCES

- Burak Bartan and Mert Pilanci. Neural spectrahedra and semidefinite lifts: Global convex optimization of polynomial activation neural networks in fully polynomial-time. arXiv preprint arXiv:2101.02429, 2021a.
- Burak Bartan and Mert Pilanci. Training quantized neural networks to global optimality via semidefinite programming. International Conference on Machine Learning (ICML), 2021b.
- Lénaïc Chizat and Francis Bach. On the global convergence of gradient descent for over-parameterized models using optimal transport. Advances in Neural Information Processing Systems, 31:3036–3046, 2018.
- Thomas M Cover. Geometrical and statistical properties of systems of linear inequalities with applications in pattern recognition. IEEE transactions on electronic computers, 14(3):326–334, 1965.
- Tolga Ergen and Mert Pilanci. Implicit convex regularizers of cnn architectures: Convex optimization of two-and three-layer networks in polynomial time. International Conference on Learning Representations (ICLR), 2020.
- Tolga Ergen and Mert Pilanci. Convex geometry and duality of over-parameterized neural networks. Journal of Machine Learning Research, 22(12):1–63, 2021a.
- Tolga Ergen and Mert Pilanci. Global optimality beyond two layers: Training deep relu networks via convex programs. International Conference on Machine Learning, pp. 2993–3003. PMLR, 2021b.
- Tolga Ergen and Mert Pilanci. Revealing the structure of deep neural networks via convex duality. In International Conference on Machine Learning, pp. 3004–3014. PMLR, 2021c.
- Tolga Ergen, Arda Sahiner, Batu Ozturkler, John Pauly, Morteza Mardani, and Mert Pilanci. Demystifying batch normalization in relu networks: Equivalent convex optimization models and implicit regularization. arXiv preprint arXiv:2103.01499, 2021.
- Suriya Gunasekar, Jason Lee, Daniel Soudry, and Nathan Srebro. Characterizing implicit bias in terms of optimization geometry. International Conference on Machine Learning, pp. 1832–1841. PMLR, 2018.
- Vikul Gupta, Burak Bartan, Tolga Ergen, and Mert Pilanci. Exact and relaxed convex formulations for shallow neural autoregressive models. International Conference on Acoustics, Speech, and Signal Processing, 2021.
- Yuanzhi Li, Tengyu Ma, and Hongyang R Zhang. Learning over-parametrized two-layer neural networks beyond ntk. Conference on Learning Theory, pp. 2613–2682. PMLR, 2020.
- Kaifeng Lyu and Jian Li. Gradient descent maximizes the margin of homogeneous neural networks. arXiv preprint arXiv:1906.05890, 2019.
- Chao Ma, Lei Wu, and Weinan E. The quenching-activation behavior of the gradient descent dynamics for two-layer neural network models. arXiv preprint arXiv:2006.14450, 2020.
- Hartmut Maennel, Olivier Bousquet, and Sylvain Gelly. Gradient descent quantizes relu network features. arXiv preprint arXiv:1803.08367, 2018.
- Song Mei, Andrea Montanari, and Phan-Minh Nguyen. A mean field view of the landscape of two-layer neural networks. Proceedings of the National Academy of Sciences, 115(33):E7665–E7671, 2018.
- Mary Phuong and Christoph H Lampert. The inductive bias of relu networks on orthogonally separable data. 2021.
- Mert Pilanci and Tolga Ergen. Neural networks are convex regularizers: Exact polynomial-time convex optimization formulations for two-layer networks. pp. 7695–7705, 2020.

Arda Sahiner, Tolga Ergen, Batu Ozturkler, Burak Bartan, John Pauly, Morteza Mardani, and Mert Pilanci. Hidden convexity of wasserstein gans: Interpretable generative models with closed-form solutions. arXiv preprint arXiv:2107.05680, 2021a.

Arda Sahiner, Tolga Ergen, John Pauly, and Mert Pilanci. Vector-output relu neural network problems are copositive programs: Convex analysis of two layer networks and polynomial-time algorithms. International Conference on Learning Representations (ICLR), 2021b.

Arda Sahiner, Morteza Mardani, Batu Ozturkler, Mert Pilanci, and John Pauly. Convex regularization behind neural reconstruction. International Conference on Learning Representations (ICLR) 2021c.

Daniel Soudry, Elad Hoffer, Mor Shpigel Nacson, Suriya Gunasekar, and Nathan Srebro. The implicit bias of gradient descent on separable data. The Journal of Machine Learning Research, 19(1): 2822–2878, 2018.

Gang Wang, Georgios B Giannakis, and Jie Chen. Learning relu networks on linearly separable data: Algorithm, optimality, and generalization. IEEE Transactions on Signal Processing, 67(9): 2357–2370, 2019.

Qiuling Yang, Alireza Sadeghi, Gang Wang, and Jian Sun. Learning two-layer relu networks is nearly as easy as learning linear classifiers on separable data. IEEE Transactions on Signal Processing, 69:4416–4427, 2021.

Mo Zhou, Rong Ge, and Chi Jin. A local convergence theory for mildly over-parameterized two-layer neural network. arXiv preprint arXiv:2102.02410, 2021.

## A DEFINITIONS AND NOTIONS

We introduce several useful definitions and notions which will be utilized in the proof.

### A.1 DEFINITIONS

**Definition 1** Let  $O \subseteq \mathbb{R}^n$  be an open set and let  $f: O \rightarrow \mathbb{R}$  be locally Lipschitz continuous at  $x \in O$ . Let  $D_F$  be the differentiable points of  $f$  in  $O$ . The B-subdifferential of  $f$  at  $x$  is defined by

$$\partial_B f(x) := \lim_{k \rightarrow 1} F^0(x^k); x^k \in D_F; \|x^k - x\|_2 \rightarrow 0 \quad (28)$$

The set  $\partial f(x) = \text{co}(\partial_B f(x))$  is called Clarke's subdifferential, where  $\text{co}$  denotes the convex hull.

**Definition 2** A matrix  $A$  is spike-free if and only if the following conditions hold: for  $\forall u \in \mathbb{R}^n, \|u\|_2 = 1$ , there exists  $\lambda \in \mathbb{R}, \lambda \geq 1$  such that

$$(Au)_+ = \lambda z \quad (29)$$

This is equivalent to say that

$$\max_{u: \|u\|_2 = 1; (I - XX^T)(Xu)_+ = 0} \|X^T(Xu)_+\|_2 \leq 1 \quad (30)$$

### A.2 NOTIONS

We use the following letters for indexing.

- The index  $n$  is for the  $n$ -th data sample  $x_n$ .
- We use the index  $i$  to represent the  $i$ -th neuron-pair  $(w_{1,i}; w_{2,i})$ .
- The index  $j$  is for the  $j$ -th masking matrix  $D_j \in \mathbb{R}^{2 \times P}$ .

## B PROOFS IN SECTION 3

### B.1 PROOF FOR PROPOSITION 2

Consider the following loss function:  $\mathbb{R}^N \times \mathbb{R}^N \rightarrow \mathbb{R}, [f + 1g]$

$$\gamma(z; y) = \begin{cases} 0; & y_n z_n \leq -1; \forall n \in [N]; \\ +1; & \text{otherwise} \end{cases} \quad (31)$$

For a given  $y \in \mathbb{R}^N, \gamma(z; y)$  is a convex loss function of  $z$ . The non-convex max-margin is equivalent to

$$\min_{W} \Gamma((XW)_+, w_2; y) + \frac{1}{2} \|W\|_F^2 + \|w_2\|_2^2 \quad (32)$$

According to Appendix A.13 in (Pilanci & Ergen, 2020), the problem (32) is equivalent to

$$\begin{aligned} \min_{\{u_i^0\}_{i=1}^P} & \Gamma\left(\sum_{i=1}^P D_i X(u_i^0 - u_i); y\right) + \frac{1}{2} \|W\|_F^2 + \|w_2\|_2^2; \\ \text{s.t. } & (2D_i - I)Xu_i \leq 0; (2D_i + I)Xu_i^0 \leq 0; \forall i \in [P]; \end{aligned} \quad (33)$$

This is equivalent to (14). For a fixed  $y \in \mathbb{R}^N, \gamma(z; y)$  with respect to  $z$  can be computed by

$$\begin{aligned} I(\hat{\cdot}; y) &= \max_{z \in \mathbb{R}^N} z^T x - \gamma(z; y) \\ &= \max_{z \in \mathbb{R}^N} z^T \hat{\cdot}; \text{ s.t. } \text{diag}(y)z \leq 1; \\ &= \begin{cases} y^T \hat{\cdot}; & \text{diag}(y) \hat{\cdot} \leq 0 \\ +1; & \text{otherwise} \end{cases} \end{aligned} \quad (34)$$

According to Theorem 6 in (Pilanci & Ergen, 2020), the dual problem of (14) writes

$$\max_{\gamma} \Gamma(\gamma); \text{ s.t. } \max_{u: \|u\|_2=1} j^T (Xu)_+ \leq 1; \quad (35)$$

which is equivalent to

$$\max_{\gamma} \gamma^T; \text{ s.t. } \text{diag}(\gamma) \succeq 0; \max_{u: \|u\|_2=1} j^T (Xu)_+ \leq 1; \quad (36)$$

By taking  $\gamma = \hat{\gamma}$ , we derive (15). This completes the proof.

## B.2 PROOF FOR PROPOSITION 3

For the non-convex max-margin problem (13), consider the Lagrange function

$$L(W_1; w_2; \gamma) = \frac{1}{2} (k W_1 k_F^2 + k w_2 k_2^2) - \gamma^T (Y (X W_1)_+ + w_2 - 1)$$

where  $\gamma \succeq 0$ . The KKT point of the non-convex max-margin problem (13) (in terms of B-subdifferential) satisfies

$$\begin{aligned} 0 &\in \partial_{W_1} L(W_1; w_2; \gamma); \\ w_2 &- (X W_1)_+^T = 0; \\ \sum_n (y_n (x_n^T W_1)_+ + w_2 - 1) &= 0; \end{aligned} \quad (37)$$

The KKT condition on the  $i$ -th column of  $W_1$  is equivalent to

$$w_{1;i} = \sum_{n=1}^N w_{2;i} x_n g_{n;i}; \quad (38)$$

where  $g_{n;i} \in \partial (z)_+ |_{z=x_n^T w_{1;i}}$ . In other words, we have

$$g_{n;i} = \begin{cases} 1(x_n^T w_{1;i} - 0); & \text{if } x_n^T w_{1;i} \in 0; \\ 2 \in [0; 1]g; & \text{if } x_n^T w_{1;i} = 0; \end{cases} \quad (39)$$

Let  $\hat{D}_i = \text{diag}([g_{1;i}; \dots; g_{N;i}])$ . Then, we can write that

$$\begin{aligned} w_{1;i} &= \sum_{n=1}^N g_{n;i} x_n w_{2;i} \\ &= w_{2;i} X^T \hat{D}_i; \end{aligned} \quad (40)$$

From the definition of  $g_{n;i}$ , we have

$$g_{n;i} x_n^T w_{1;i} = 0; \quad (41)$$

Therefore, we can compute that

$$\begin{aligned} w_{2;i} &= (X w_{1;i})_+^T \\ &= \sum_{n=1}^N 1(x_n^T w_{1;i} - 0) x_n^T w_{1;i} \\ &= \sum_{n=1}^N g_{n;i} x_n^T w_{1;i} \\ &= w_{1;i}^T X^T \hat{D}_i; \end{aligned} \quad (42)$$

In summary, we have

$$w_{1;i} = w_{2;i} X^T \hat{D}_i; \quad w_{2;i} = w_{1;i}^T X^T \hat{D}_i; \quad (43)$$

Suppose that  $w_{2;i} \in 0$ . This implies that

$$\frac{w_{1;i}}{w_{2;i}} = X^T \hat{D}_i; \quad k X^T \hat{D}_i k_2 = 1; \quad (44)$$

This completes the proof.

## B.3 PROOF FOR THEOREM 3

PROOF We can write the Lagrange function for the convex max-margin problem (14) as

$$\begin{aligned}
& L(f u_j g_{j=1}^p; f u_j^0 g_{j=1}^p; ; f z_j g_{j=1}^p; f z_j^0 g_{j=1}^p) \\
& = \sum_{j=1}^{\mathcal{X}^p} (k u_j k_2 + k u_j^0 k_2) + \sum_{j=1}^{\mathcal{X}^p} \text{diag}(y) @ 1 \text{diag}(y) D_j X (u_j - u_j^0) A \\
& \quad + \sum_{j=1}^{\mathcal{X}^p} (z_j^T (2D_j - I) X u_j + (z_j^0)^T (2D_j - I) X u_j^0) \\
& = \sum_{j=1}^{\mathcal{X}^p} \text{diag}(y) + \sum_{j=1}^{\mathcal{X}^p} (k u_j k_2 + k u_j^0 k_2) + \sum_{j=1}^{\mathcal{X}^p} (u_j^0)^T (X^T D_j - X^T (2D_j - I) z_j^0) \\
& \quad + \sum_{j=1}^{\mathcal{X}^p} (u_j)^T (X^T D_j - X^T (2D_j - I) z_j):
\end{aligned} \tag{45}$$

where  $z_j; z_j^0 \in \mathbb{R}^N$  satisfies that  $z_j \geq 0, z_j^0 \geq 0$  for  $j \in [p]$  and  $\sum_{j=1}^{\mathcal{X}^p} \text{diag}(y) = 0$ . The KKT point shall satisfy the following KKT conditions:

$$\begin{aligned}
& X^T D_j + X^T (2D_j - I) z_j^0 \leq @_j k u_j^0 k_2; \\
& X^T D_j + X^T (2D_j - I) z_j \leq @_j k u_j k_2; \\
& \sum_{j=1}^{\mathcal{X}^p} (D_j)_{n,n} X_n^T (u_j - u_j^0) - y_n A = 0; \\
& z_{j,n} (2(D_j)_{n,n} - 1) X_n^T u_j = 0; \\
& z_{j,n}^0 (2(D_j)_{n,n} - 1) X_n^T u_j^0 = 0;
\end{aligned} \tag{46}$$

Let  $(w_{1;}; w_{2;})$  be the KKT point of the non-convex problem (2) and satisfies (17). Let  $\hat{D}_i$  be the diagonal matrix defined in Proposition 3 with respect to  $w_{1;}$  and denote  $\hat{P} = \{i \in [m] \mid w_{1; i} > w_{2; i}\}$ . Without the loss of generality, we may assume that  $w_{1; i}$  are different. (Otherwise, we can merge two neurons  $w_{1; i_1}$  and  $w_{1; i_2}$  with  $D_{i_1} = D_{i_2}$  together.)

Suppose that  $D_j \in \hat{P}$ , i.e.,  $D_j = \hat{D}_i$  for certain  $i \in [m]$ . By letting  $u_j^0 = w_{1; i} w_{2; i}, z_j^0 = 0, u_j = w_{1; i} w_{2; i}$  and  $z_j = 0$ , the following identities hold.

$$X^T D_j + X^T (2D_j - I) z_j^0 = X^T \hat{D}_i = \frac{w_{1; i}}{w_{2; i}} = \frac{u_i^0}{k u_i^0 k}; \tag{47}$$

$$X^T D_j + X^T (2D_j - I) z_j = X^T \hat{D}_i = \frac{w_{1; i}}{w_{2; i}} = \frac{u_i}{k u_i k}; \tag{48}$$

Therefore, for index  $j$  satisfying  $D_j \in \hat{P}$ , the first two KKT conditions in (46) hold.

For  $D_j \notin \hat{P}$ , we can let  $u_j = u_j^0 = 0$ . As satisfies (17), we have

$$\max_{k u_k k_2 \leq 1; (2D_j - I) X u_j^0} j^T D_j X u_j \leq 1; \tag{49}$$

According to Lemma 4 in (Pilanci & Ergen, 2020), this implies that there exist  $\epsilon > 0$  such that

$$k X^T D_j + Z^T (2D_j - I) z_j^0 k \leq 1; k X^T D_j + Z^T (2D_j - I) z_j k \leq 1; \tag{50}$$

Therefore, the first two KKT conditions in (46) hold.

From our choice of  $u_j; z_j; u_j^0; z_j^0$ , the last two KKT conditions in (46) hold. We also note that

$$\sum_{j=1}^{\mathcal{X}^p} D_j X (u_j^0 - u_j) = \sum_{i=1}^{\mathcal{X}^n} (X w_{1; i})_+ w_{2; i}; \tag{51}$$

As  $(w_{1;}; w_{2;})$  is the KKT point of the non-convex problem, the third KKT condition (46) holds. This completes the proof.

## C PROOFS IN SECTION 3.1

In this section, we present several proofs for propositions in Section 3.1.

### C.1 PROOF FOR PROPOSITION 4

We start with two lemmas.

**Lemma 5** Suppose that  $u_0 = X^T \hat{D}_0$  and  $\|u_0\|_2 = 1$ . For any masking matrix  $D_j \in \mathbb{P}$  such that  $(D_j - \hat{D}_0)(I(\cdot > 0)) = 0$ , we have

$$\max_{(2D_j - I)Xu_0; \|u\|_2 = 1} X^T D_j Xu = 1: \quad (52)$$

**PROOF** According to Lemma 4 in (Pilanci & Ergen, 2020), the constraint (52) is equivalent to that there exist  $z_j \in \mathbb{R}^N$  such that  $z_j \geq 0$  and

$$\|X^T D_j + X^T (2D_j - I)z_j\|_2 = 1: \quad (53)$$

Consider the index  $n \in [N]$  such that  $(D_j - \hat{D}_0)_{nn} \neq 0$ . As  $(D_j - \hat{D}_0)(I(\cdot > 0)) = 0$ , we have  $(D_j - \hat{D}_0)_{nn} \leq 0$ . We let  $(z_j)_n = \frac{1}{(D_j - \hat{D}_0)_{nn}}$ . If  $(\hat{D}_0)_{nn} = 0$ , then we have  $(D_j)_{nn} = 1$  and

$$(D_j - \hat{D}_0)_{nn} (z_j)_n = 0 = X_n^T (2(D_j)_{nn} - 1)(z_j)_n: \quad (54)$$

If  $(\hat{D}_0)_{nn} = 1$ , then we have  $(D_j)_{nn} = 0$  and

$$(D_j - \hat{D}_0)_{nn} (z_j)_n = 0 = X_n^T (2(D_j)_{nn} - 1)(z_j)_n: \quad (55)$$

For other index  $n \in [N]$ , we simply let  $(z_j)_n = 0$ . Then, we have

$$(D_j - \hat{D}_0)_{nn} (z_j)_n = 0 = X_n^T (2(D_j)_{nn} - 1)(z_j)_n: \quad (56)$$

Based on our choice of  $z_j$ , we have  $z_j \geq 0$  and for  $n \in [N]$

$$(D_j - \hat{D}_0)_{nn} (z_j)_n = X_n^T (2(D_j)_{nn} - 1)(z_j)_n: \quad (57)$$

This implies that

$$X^T (D_j - \hat{D}_0) = X^T (2D_j - I)z_j: \quad (58)$$

Hence, we have

$$X^T D_j + X^T (2D_j - I)z_j = X^T \hat{D}_0 = u_0: \quad (59)$$

Therefore  $\|X^T D_j + X^T (2D_j - I)z_j\|_2 = 1$ .

**Lemma 6** Suppose that the data is orthogonal separable and  $u_0 = X^T \hat{D}_0$  and  $\|u_0\|_2 = 1$ . For any masking matrix  $D_j$  such that  $(D_j - \hat{D}_0)(I(\cdot > 0)) = 0$ , we have  $\|X^T D_j\|_2 \leq \|u_0\|_2 = 1$ . Therefore, (52) holds.

**PROOF** We note that  $u_0 = X^T (\hat{D}_0 - D_j) + X D_j$ . Denote  $a = X^T (\hat{D}_0 - D_j)$  and  $b = X^T D_j$ . We note that

$$a^T b = \sum_{n: (\hat{D}_0)_{nn} = 1; (D_j)_{nn} = 0} X_n^T a_n + \sum_{n: (\hat{D}_0)_{nn} = 0; (D_j)_{nn} = 0} X_n^T a_n + \sum_{n: (\hat{D}_0)_{nn} = 0; (D_j)_{nn} = 1} X_n^T a_n + \sum_{n: (\hat{D}_0)_{nn} = 1; (D_j)_{nn} = 1} X_n^T a_n = \sum_{n: (\hat{D}_0)_{nn} = 1; (D_j)_{nn} = 1} X_n^T a_n = 0: \quad (60)$$

As  $\text{diag}(y) \geq 0$ ,  $\sum_n y_n$  has the same signature with  $\sum_n y_n^2$ . Therefore, from the orthogonal separability of the data, we have

$$\sum_n X_n^T a_n X_n^T a_n = 0: \quad (61)$$

This immediately implies that  $a^T b = 0$ . Therefore,

$$\|u_0\|_2^2 = \|a + b\|_2^2 = \|a\|_2^2 + \|b\|_2^2 + 2a^T b = \|a\|_2^2 + \|b\|_2^2: \quad (62)$$

This completes the proof.

Based on Lemma 5 and Lemma 6, we present the proof for Proposition 3. Let  $w_{1,+} = \frac{w_{1,+}}{w_{2,+}}$ . From the proof of Proposition 3, we note that  $\|u_+\|_2 = 1$ . For any masking matrix  $D_j \in \mathbb{P}$ , let  $\hat{D} = \hat{D}_{i,+} D_j$ . As  $\hat{D}_{i,+} \in \mathbb{P}$ , according to Lemma 6, we have

$$\|X^T \hat{D}\|_2 \leq \|X^T \hat{D}_{i,+}\|_2 = \|u_+\|_2 = 1: \quad (63)$$

As  $Y \geq 0$  and  $\hat{D}_{i,+} = \text{diag}(I(y = 1))$ , we have  $(D_j - \hat{D})(I(\cdot > 0)) = D_j (I(\hat{D}_{i,+})) (I(\cdot > 0)) = 0$ . From Lemma 5, we note that  $\hat{D}$  satisfies (52). Similarly, we can show that  $\hat{D}_{i,+}$  also satisfies (52). This completes the proof.

## C.2 PROOF FOR PROPOSITION 5

PROOF Note that  $Y \succeq 0$ . Let  $Y_+ = \text{diag}(I(y = 1))$  and  $Y_- = \text{diag}(I(y = -1))$ . We claim that

$$\max_{\|u\|_2 = 1} \langle Xu \rangle_+ = \max_{\|u\|_2 = 1} \langle Y_+ Xu \rangle_+ \quad (64)$$

Firstly, we note that

$$\langle Xu \rangle_+ = \sum_{n=1}^N (x_n^T u)_+ = \sum_{n=1}^N (x_n)_+^T (x_n^T u)_+ = \langle Y_+ Xu \rangle_+ \quad (65)$$

This implies that  $\max_{\|u\|_2 = 1} \langle Xu \rangle_+ = \max_{\|u\|_2 = 1} \langle Y_+ Xu \rangle_+$ .

On the other hand, suppose that  $z = \arg \max_{\|u\|_2 = 1} \langle Y_+ Xu \rangle_+$ . As  $X$  is spike-free, there exists  $z$  such that  $\|z\|_2 = 1$  and  $Xz = (Xu)_+$ . Therefore, we have

$$\langle Y_+ Xu \rangle_+ = \langle Y_+ Xz \rangle_+ = \langle Xz \rangle_+ = \langle (Xz)_+ \rangle_+ \quad (66)$$

This implies that  $\max_{\|u\|_2 = 1} \langle Xu \rangle_+ = \max_{\|u\|_2 = 1} \langle Y_+ Xu \rangle_+$ .

For any  $D_j \succeq P$  with  $D_j \succeq Y_+$ . We note that

$$\langle Xu \rangle_+ \leq \langle D_j Xu \rangle_+ \leq \langle Y_+ Xu \rangle_+ \quad (67)$$

Combining with (65), this implies that  $\max_{\|u\|_2 = 1} \langle Xu \rangle_+ = \max_{\|u\|_2 = 1} \langle D Xu \rangle_+$ .

Let us go back to the original problem. Let  $\hat{D}_{i_+} = \frac{w_{1;i_+}}{w_{2;i_+}}$ . We note that  $\langle Xw_+ \rangle_+ = \hat{D}_{i_+} \langle Xw_+ \rangle_+ = \hat{D}_{i_+} \langle XX^T \hat{D}_{i_+} w_+ \rangle_+$ . Therefore, we have

$$\langle Xw_+ \rangle_+ = \langle \hat{D}_{i_+} XX^T \hat{D}_{i_+} w_+ \rangle_+ = \langle kX^T \hat{D}_{i_+} k^2 w_+ \rangle_+ = \langle ku_+ k^2 w_+ \rangle_+ = 1 \quad (68)$$

Thus, for any  $\|u\|_2 = 1$ , suppose that  $\langle Xu \rangle_+ = \langle Xz \rangle_+$ , where  $\|z\|_2 = 1$ . Then, we have

$$\langle \hat{D}_{i_+} Xu \rangle_+ = \langle \hat{D}_{i_+} Xz \rangle_+ \leq \langle kzk_2 \rangle_+ = 1 \quad (69)$$

Therefore,  $\max_{\|u\|_2 = 1} \langle Xu \rangle_+ = \max_{\|u\|_2 = 1} \langle D_+ Xu \rangle_+ = 1$ . Similarly, we have

$$\min_{\|u\|_2 = 1} \langle Xu \rangle_+ = \min_{\|u\|_2 = 1} \langle D_- Xu \rangle_+ = -1$$

This completes the proof.

## D PROOFS IN SECTION 4

### D.1 PROOF FOR LEMMA 1

PROOF According to the sub-gradient flow (22), we can compute that

$$\frac{\partial}{\partial t} \langle kw_{1;i} k^2 - w_{2;i}^2 \rangle = 2w_{1;i}^T w_{2;i} g(u; e) - 2w_{2;i} w_{1;i}^T g(u; e) = 0 \quad (70)$$

Let  $T_0 = \sup\{t \mid \langle kw_{1;i}(t) k^2 - w_{2;i}(t)^2 \rangle > 0; i \in [n]; t \in [0; T]\}$ . For  $t \in [0; T_0)$ , as the neural network is scaled, it is sufficient to study the dynamics of  $w_{1;i}$  in the polar coordinate. Let us write  $w_{1;i}(t) = e^{f_i(t)} u_i(t)$ , where  $\|u_i(t)\|_2 = 1$ . Then, in terms of polar coordinate, the projected gradient flow follows

$$\begin{aligned} \frac{\partial}{\partial t} f_i &= \text{sign}(w_{2;i}) u_i^T g(u_i; e); \\ \frac{\partial}{\partial t} u_i &= \text{sign}(w_{2;i}) g(u_i; e) - u_i^T g(u_i; e) u_i. \end{aligned} \quad (71)$$

Without the loss of generality, we may assume that  $w_{2;i}(0) \leq 0$  for  $i \in [m]$ . Denote

$$x_{\max} = \max_{i \in [n]} \|x_i\|_2 \quad (72)$$



From the definition of  $e$ , we have  $\|e\|_2 = 1$ . Therefore, we have

$$\frac{\partial}{\partial t} \|g(u_i; e)\|_2 = \sum_{j: x_j^T u > 0} \tilde{x}_j x_j \frac{nx_{\max}}{2} \quad (73)$$

Therefore, for  $\forall t \in \mathbb{R}^+$ , we have

$$r_i(t) = r_i(0) + \frac{nx_{\max} t}{4} \quad (74)$$

which implies that  $\forall w_{2,i}(t) > 0$ . This implies that  $\forall_0 = 1$ .

## D.2 PROOF OF LEMMA 2

PROOF As we have  $\|w_{1,i}\|_2 = \|w_{2,i}\|_2$ , for  $n \in [N]$ , we can compute that

$$\begin{aligned} \|q_n\|_2 &= \|(x_n^T W_{1,+} + w_2)\|_2 \\ &= \sum_{i=1}^n \|x_n^T w_{1,i} + w_{2,i}\|_2 \\ &= \sum_{i=1}^n \|w_{1,i}\|_2 \|x_n^T w_{1,i} + w_{2,i}\|_2 \\ &= \sum_{i=1}^n \|w_{1,i}\|_2 \|x_n^T w_{1,i}\|_2 \\ &= \sum_{i=1}^n \|x_n\|_2 \|w_{1,i}\|_2^2 \end{aligned} \quad (75)$$

Note that  $\tilde{y}_n = y_n - \eta(q_n)$  and  $\tilde{y}_n = y_n - \eta(0)$ . As  $\eta$  is  $\frac{1}{4}$ -Lipschitz continuous, we have

$$\|\tilde{y}_n - y_n\|_2 \leq \frac{1}{4} \|q_n\|_2 \leq \frac{\|x_n\|_2}{4} \sum_{i=1}^n \|w_{1,i}\|_2^2 \quad (76)$$

For any  $\epsilon \in [0, 1]$ , as  $\tilde{y}_n \in [0, 1]$  for  $n \in [N]$ , we have

$$\begin{aligned} \|g(\tilde{y}; e)\|_2 &= \|g(\tilde{y}; y)\|_2 \\ &= \sum_{k: \tilde{y}_k > 0} \tilde{x}_k x_k \\ &= \sum_{k=1}^n \frac{\|x_k\|_2^2}{4} \sum_{j=1}^n \|w_{1,j}\|_2^2 \\ &= c_1 \sum_{i=1}^n \|w_{1,i}\|_2^2 \end{aligned} \quad (77)$$

where  $c_1 = \frac{1}{4} \|x\|_2^2 > 0$  is a constant. Therefore, we can bound  $\|g(\tilde{y}; e(t))\|_2$  by

$$\|g(\tilde{y}; e(t))\|_2 \leq \|g(\tilde{y}; y)\|_2 + c_1 \sum_{i=1}^n \|w_{1,i}(t)\|_2^2 \leq d_{\max} + c_1 \sum_{i=1}^n e^{2r_i(t)} \quad (78)$$

where we let

$$d_{\max} = \max_{2Q} \|g(\tilde{y}; y)\|_2 \quad (79)$$

Let  $r(t) = \max_{i \in [m]} r_i(t)$ . We note that

$$\frac{d}{dt} r(t) \leq d_{\max} + c_1 n e^{2r(t)} - c_2 (1 + e^{2r(t)}); \quad (80)$$

where  $c_2 = \max\{nc_1, d_{\max}\} > 0$  is a constant. If we start with  $r(0) = 0$ , then,  $r(t)$  cannot grow much faster than  $2t$ . Let  $r(t)$  satisfy the following ODE:

$$\frac{d}{dt} r(t) = c_2 (1 + e^{2r(t)}); \quad (81)$$

The solution is given by

$$r_a(t) = c_2(t - a) + \frac{1}{2} \log(1 + e^{2c_2(t - a)}); \quad (82)$$

where  $a > 0$  is a parameter depending on the initialization. For any initial  $r(0)$ , we have a unique solution satisfying  $r_a(0) = r(0)$ . Therefore, we have  $r(t) \leq r_a(t)$  and

$$kg(\cdot; e(t)) \leq g(\cdot; y=4)k_2 - c_1 n e^{2r_a(t)}; \quad (83)$$

According to the bound (83), by choosing a sufficiently small  $a_0$ , (which leads to a sufficiently small  $a$ ), such that

$$e^{2r_a(T)} \leq \min\left\{\frac{d_{\min}}{16c_1}, \frac{d_{\min}^2}{4n^2x_{\max}^3}\right\}; \quad (84)$$

Therefore, for  $t \leq T$ , we have

$$kg(\cdot; e(t)) \leq g(\cdot; y=4)k_2 - c_1 n e^{2r_a(t)} \leq c_1 n e^{2r_a(T)} \leq \frac{d_{\min}}{8}; \quad (85)$$

Hence, we have

$$kg(\cdot(t); e(t)) \leq g(\cdot(t); y=4)k_2 - c_1 n e^{2r_a(t)} \leq c_1 n e^{2r_a(T)} \leq \frac{d_{\min}}{8}; \quad (86)$$

We can compute that

$$\frac{d}{dt} \tilde{r}_i = |y_i|^{(2)}(q) \frac{d}{dt} q; \quad (87)$$

As  $|y_i|^{(2)}(q) \in [0, 1]$ , we can compute that

$$\begin{aligned} \frac{d}{dt} q &= \sum_{j=1}^n |w_{2,j}| k_{x_i} k_2 \frac{d}{dt} w_{1,j} \\ &+ \sum_{j=1}^n |w_{1,j}| k_2 k_{x_i} k_2 \frac{d}{dt} w_{2,j} \\ &\leq \frac{n}{4} \sum_{j=1}^n |w_{1,j}| k_2^2 x_{\max}^2 + \frac{n}{4} \sum_{j=1}^n |w_{2,j}| x_{\max}^2 \\ &\leq \frac{nx_{\max}^2}{2} e^{2r(t)}. \end{aligned} \quad (88)$$

Therefore, we have

$$\frac{d}{dt} \tilde{r}_i \leq |y_i|^{(2)}(q) \frac{d}{dt} q \leq \frac{1}{4} \frac{d}{dt} q \leq \frac{nx_{\max}^2}{8} e^{2r(t)}; \quad (89)$$

Suppose that  $\text{sign}(X u(s)) = \text{sign}(t)$  holds for  $s$  in a small neighborhood. Then, we have

$$\begin{aligned} \frac{d}{dt} g(u(t); e(t)) &= \frac{d}{dt} g(\cdot(t); e(t)) \\ &\leq \sum_{i=1}^n |k_{x_i} k_2| \frac{d}{dt} \tilde{r}_i \leq \frac{n^2 x_{\max}^3}{8} e^{2r(t)} \\ &\leq \frac{n^2 x_{\max}^3}{8} e^{2r_a(T)} \leq \frac{d_{\min}^2}{16}. \end{aligned} \quad (90)$$

This completes the proof.

### D.3 PROOF OF LEMMA 3

PROOF Let  $T_0 = \sup\{T \mid \text{sign}(Xu(t)) = \text{sign}(Xu_0); \forall t \in [0; T]\}$ . We analyze the dynamics of  $u(t)$  in the interval  $[0; \min\{T_0; T\}]$ . For  $t \in [0; \min\{T_0; T\}]$ , as the statements in Lemma 2 hold, we can compute that

$$\begin{aligned} & \frac{d}{dt} v(t)^\top u(t) \\ &= \frac{d}{dt} \frac{g(\cdot; e(t))^\top}{kg(\cdot; e(t))k_2} u(t) \\ &= \frac{g(\cdot; e(t))^\top}{kg(\cdot; e(t))k_2} \frac{d}{dt} u(t) + u(t)^\top \frac{1}{kg(\cdot; e(t))k_2} \frac{d}{dt} g(\cdot; e(t)) \\ & \quad - u(t)^\top g(\cdot; e(t)) \frac{g(\cdot; e(t))^\top \frac{d}{dt} g(\cdot; e(t))}{kg(\cdot; e(t))k_2^3} \\ & \leq \frac{g(\cdot; e(t))^\top \frac{d}{dt} u(t)}{g_{\min}} - \frac{2}{g_{\min}} \frac{d}{dt} g(\cdot; e(t)) \\ & \leq k g(\cdot; e(t))k_2^{-1} (v(t)^\top u(t))^2 - \frac{g_{\min}}{8} \\ & \leq g_0 (1 - \frac{1}{8}) (v(t)^\top u(t))^2 - \frac{g_{\min}}{8} \\ & \leq g_0 (1 - \frac{1}{4}) (v(t)^\top u(t))^2 : \end{aligned} \tag{91}$$

Here we utilize that  $g_0 \geq g_{\min}$ , where  $g_{\min}$  is defined in (24). Let  $z(t)$  satisfies the ODE

$$\frac{dz(t)}{dt} = (1 - \frac{1}{4} z(t)^2) g_0; \tag{92}$$

with initialization  $z(0) = v_0^\top u_0$ . Then, we note that

$$z(t) = \frac{1}{1 - \frac{1}{8}} \frac{2^{\frac{p}{1 - \frac{1}{8}}}}{1 + c_3 \exp(2g_0 t (1 - \frac{1}{8}))}; \tag{93}$$

where  $c_3 = \frac{1 - \frac{1}{8}}{v_0^\top u_0} (1 - \frac{1}{8})$ . We can compute that

$$z(T_3) = 1 : \tag{94}$$

According to the comparison theorem, for  $t \in [0; T_3]$ , we have

$$v(t)^\top u(t) \leq z(t); \tag{95}$$

We first consider the case where  $T_0 = 1$ . As  $T_0 = 1$ , we have

$$v(T)^\top u(T) \leq z(T) = 1 : \tag{96}$$

Therefore, the second event holds for  $T$ .

Otherwise, we have  $T_0 < 1$ . Recall that  $u_1 = u(T_0)$  and  $v_1 = \lim_{t \rightarrow T_0} v(t)$ . Let  $T_1 = \sup\{T \mid v(t)^\top u(t) < v_1^\top u_1; \forall t \in [0; T]\}$  and  $T_2 = \sup\{T \mid v(t)^\top u(t) < 1; \forall t \in [0; T]\}$ . If  $T_2 > T_0$ , for  $t \in [0; T_2]$ , we have

$$\frac{d}{dt} v(t)^\top u(t) \geq (1 - \frac{1}{4})^2 g_0 > 0; \tag{97}$$

Therefore,  $v(t)^\top u(t)$  monotonically increases on  $[0; T_2]$ . As  $v(t)^\top u(t) \leq z(t)$  for  $t \in [0; T_0]$ , we have that  $z(T_2) = v(T_2)^\top u(T_2) = 1 = z(T_3)$ . Hence, we have  $T_2 = T$ . Therefore, the second condition of the first event holds at  $T = T_2$ .

Then, we consider the case where  $T_0 < 1$ . For  $t \in [0; T_0]$ , we have  $v(t)^\top u(t) \leq 1$ . This implies that  $v_1^\top u_1 \leq 1$ . Apparently, we have  $T_1 = T_0$ . If  $T_1 < T_0$ , as  $T_0 = T_2$ , for  $t \in [0; T_0]$ , the inequality

(97) holds. This implies that  $\lim_{t \rightarrow T_0} v(t)^T u(T_0) > v(T_1)^T u(T_1) = \lim_{t \rightarrow T_0} v(t)^T u(T_0)$ , which leads to a contradiction. Therefore, we have  $T_1 = T_0$ . We note that

$$z(T^{\text{shift}}(u_1^T v_1)) = u_1^T v_1: \quad (98)$$

As  $u(t)^T g(u(t); e(t)) = z(t)$  for  $t \in [0; T_0]$ , we have that  $z(T_1) = u_1^T v_1 = z(T_0)$ . Hence, we have  $T_0 = T_1 = T^{\text{shift}}(u_1^T v_1)$ . This completes the proof.

#### D.4 PROOF OF PROPOSITION 6

We first introduce a lemma.

**Lemma 7** Let  $a, b \in \mathbb{R}^d$  and  $0 < c < \|a\|$ . Suppose that  $\|b\|_2 \leq c$  and  $\|a\|_2 \geq c$ . Then, we have

$$\frac{a}{\|a\|_2} \cdot \frac{b}{\|b\|_2} \geq \frac{2}{c}: \quad (99)$$

**PROOF** As  $\|a\|_2 \geq c$ , we have  $\|b\|_2 \leq \|a\|_2 \cdot \frac{c}{\|a\|_2} < \|a\|_2$ . We first note that

$$\|a\|_2^{-1} \|b\|_2^{-1} = \frac{\|a\|_2 \|b\|_2}{\|a\|_2 \|b\|_2} \geq \frac{c}{\|a\|_2 \|b\|_2}: \quad (100)$$

Therefore, we can compute that

$$\begin{aligned} & \frac{a}{\|a\|_2} \cdot \frac{b}{\|b\|_2} \geq \frac{c}{\|a\|_2 \|b\|_2} \\ & \frac{a}{\|a\|_2} \cdot \frac{b}{\|a\|_2} + \frac{1}{\|a\|_2} \cdot \frac{1}{\|b\|_2} \|b\|_2 \\ & \frac{c}{c} + \frac{c}{c} = \frac{2}{c}: \end{aligned} \quad (101)$$

This completes the proof.

Then we present the proof of Proposition 6.

**PROOF** As  $y^T (Xu_0)_+ > 0$ , with sufficiently small initialization and sufficiently small  $\epsilon > 0$ , we also have  $e(0)^T (Xu_0)_+ = y^T (Xu_0)_+ - \epsilon \|X\|_2 k_2 e(0) > 0$ . We prove that there exists a time  $T$  such that  $u(T)^T v(T) \geq \frac{3}{4}$  by contradiction. Denote  $v_0 = v(0)$ . For all possible values of  $kg(u; y)k_2$ , we can arrange them from the smallest to the largest by  $g^{(1)} < \dots < g^{(p)}$ .

Let  $T_i = \frac{p-1}{2} \frac{1}{1-\epsilon g^{(i)}} \log p \frac{1-\epsilon g^{(i+1)}}{1-\epsilon g^{(i)}} = 2 \log p \frac{1-\epsilon g^{(i+1)}}{1-\epsilon g^{(i)}} \frac{1}{1-\epsilon g^{(i)}} v_0^T u_0$  and  $T = \sum_{i=1}^p T_i$ . Suppose that  $r_0$  is sufficiently small such that statements in Lemma 2 holds. According to Lemma 3, we can find  $0 = t_0 < t_1 < \dots$  such that for  $i = 1, \dots, p$ ,  $\text{sign}(Xu(t))$  is constant on  $(t_{i-1}; t_i)$  and  $\text{sign}(Xu(t_{i-1})) \neq \text{sign}(Xu(t_i))$ . We write  $u_i = u(t_i)$ ,  $g_i = kg(u(t_i); y)k_2$ ,

$$g_i = \lim_{t \rightarrow t_i} g(u(t); e(t)); \quad \tilde{g}_i = g(u(t_i); e(t)); \quad (102)$$

$v_i = \frac{g_i}{kg_i k_2}$  and  $\tilde{v}_i = \frac{\tilde{g}_i}{kg_i k_2}$ . We note that  $\tilde{g}_i = g_i^+$ . According to Lemma 3, we have

$$\begin{aligned} t_i - t_{i-1} & \geq \frac{1}{2} \frac{1}{1-\epsilon g_{i-1}} \log p \frac{1-\epsilon g_i}{1-\epsilon g_{i-1}} \frac{1}{1-\epsilon g_{i-1}} (v_i)^T u_i - \frac{1}{2} \frac{1}{1-\epsilon g_{i-1}} \log p \frac{1-\epsilon g_i}{1-\epsilon g_{i-1}} \frac{1}{1-\epsilon g_{i-1}} (v_{i-1})^T u_{i-1} \\ & \geq \frac{1}{2} \frac{1}{1-\epsilon g_{\min}} \log p \frac{1-\epsilon g_i}{1-\epsilon g_{i-1}} \frac{1}{1-\epsilon g_{i-1}} (v_i)^T u_i - \frac{1}{2} \frac{1}{1-\epsilon g_{\min}} \log p \frac{1-\epsilon g_i}{1-\epsilon g_{i-1}} \frac{1}{1-\epsilon g_{i-1}} (v_{i-1})^T u_{i-1} : \end{aligned} \quad (103)$$

Here we utilize that  $g_{i-1} \geq g_{\min}$ , where  $g_{\min}$  is defined in (24). This implies that

$$\frac{1}{1-\epsilon g_{i-1}} \frac{1}{1-\epsilon g_{i-1}} (v_i)^T u_i \geq e^{2 \frac{1-\epsilon g_{\min}}{1-\epsilon g_{i-1}} (t_i - t_{i-1})} \frac{1}{1-\epsilon g_{i-1}} \frac{1}{1-\epsilon g_{i-1}} (v_{i-1})^T u_{i-1}. \quad (104)$$

We can show that for satisfying  $\log \frac{1 + g_{\min}^{-1} v_0^T u_0}{1 - g_{\min}^{-1} v_0^T u_0} \geq 2$  and  $t \in T$ , we have  $kg(u(t); y=4)k_2 > g_{\min}$ . According to Lemma 3, as  $g_{\min}$ , we have

$$\frac{1 + g_{\min}^{-1} (g_i)^T u_i}{1 - g_{\min}^{-1} (g_i)^T u_i} \leq e^{2g_{\min}^{-1} (g_i)^T u_i} \frac{1 + g_{\min}^{-1} v_0^T u_0}{1 - g_{\min}^{-1} v_0^T u_0}. \quad (105)$$

This implies that

$$\frac{1 + g_{\min}^{-1} (g_i)^T u_i}{1 - g_{\min}^{-1} (g_i)^T u_i} \leq e^{2g_{\min}^{-1} (g_i)^T u_i} \frac{1 + g_{\min}^{-1} v_0^T u_0}{1 - g_{\min}^{-1} v_0^T u_0}, \quad (106)$$

or equivalently, for any  $\psi > 0$ , we have

$$\frac{1 + g_{\min}^{-1} (g(u(t); e(t)))^T u(t)}{1 - g_{\min}^{-1} (g(u(t); e(t)))^T u(t)} \leq e^{2\psi} \frac{1 + g_{\min}^{-1} v_0^T u_0}{1 - g_{\min}^{-1} v_0^T u_0}. \quad (107)$$

Here we utilize that  $(g(u(t); e(t)))^T u(t)$  is continuous w.r.t.  $t$ . Therefore, for  $\frac{1}{2g_{\min}} \log \frac{1 + g_{\min}^{-1} v_0^T u_0}{1 - g_{\min}^{-1} v_0^T u_0} \geq \psi$ , we have

$$\frac{1 + g_{\min}^{-1} (g(u(t); e(t)))^T u(t)}{1 - g_{\min}^{-1} (g(u(t); e(t)))^T u(t)} \leq \frac{1 + g_{\min}^{-1} v_0^T u_0}{1 - g_{\min}^{-1} v_0^T u_0}. \quad (108)$$

This implies that

$$g_{\min}^{-1} (g(u(t); e(t)))^T u(t) \leq \frac{1 + g_{\min}^{-1} v_0^T u_0}{1 - g_{\min}^{-1} v_0^T u_0}. \quad (109)$$

If  $kg(u(t); y=4)k_2 = g_{\min}$ , as the statements in Lemma 2 hold, we can compute that

$$kg(u(t); y=4)k_2 = g(u(t); e(t))k_2 \frac{g_{\min}}{4} = \frac{1}{4}kg(u(t); y=4)k_2; \quad (110)$$

which implies that

$$kg(u(t); e(t))k_2 = (1 + \frac{1}{4})kg(u(t); y=4)k_2; \quad (111)$$

Therefore, we have

$$\begin{aligned} v(t)^T u(t) &= \frac{(g(u(t); e(t)))^T u(t)}{kg(u(t); e(t))k_2} \\ &= \frac{1}{1 + \frac{1}{4}} \frac{(g(u(t); e(t)))^T u(t)}{g_{\min}} \\ &= \frac{1}{1 + \frac{1}{4}} \frac{1}{g_{\min}} \frac{3}{4}. \end{aligned} \quad (112)$$

This leads to a contradiction.

Analogously, we can show that for  $\bigcup_{i=1}^p T_i$ , we have  $kg(u(t); y=4)k_2 > g_{(i)}$ . Thus, by taking  $t \in \bigcup_{i=1}^p T_i$ , we have  $kg(u(t); y=4)k_2 > g_{(p)} = g_{\max}$ . However, from the definition of  $g_{\max}$ , we have  $kg(u(t); y=4)k_2 \leq g_{\max}$ . This leads to a contradiction. Therefore, there exists a time  $T = \bigcup_{i=1}^p T_i = O(\log^{-1})$  such that  $v(T)^T u(T) \leq \frac{3}{4}$ .

We note that  $kg(u(T); y=4)k_2 = g_{\min}$ . As the statements in Lemma (2) hold, we have

$$kg(u(T); y=4)k_2 = g(u(T); e(T))k_2 \frac{g_{\min}}{8} \quad (113)$$

According to Lemma 7, we have

$$\frac{kg(u(T); y=4)k_2}{kg(u(T); y=4)k_2} = \frac{g(u(T); e(T))}{kg(u(T); e(T))k_2} \frac{2kg(u(T); y=4)k_2}{g_{\min}} \leq \frac{3}{4}. \quad (114)$$

This implies that

$$\begin{aligned} & u(T)^T \frac{g(u(T); y=4)}{kg(u(T); y=4)k_2} \\ & u(T)^T v(T) = \frac{g(u(T); y)k_2}{kg(u(T); y)k_2} = \frac{g(u(T); e(T))}{kg(u(T); e(T))k_2} \quad (115) \\ & 1 : \end{aligned}$$

Hence, we have

$$\cos \angle(u(T); g(u(T); y)) = u(T)^T \frac{g(u(T); y)}{kg(u(T); y)k_2} = u(T)^T \frac{g(u(T); y=4)}{kg(u(T); y=4)k_2} \quad 1 :$$

This completes the proof.

#### D.5 PROOF OF LEMMA 4

PROOF This is proved in Lemma 2 in (Phuong & Lampert, 2021). Here we provide an alternative proof. It is sufficient to prove for the case of local maximizer. Suppose  $w$  is a local maximizer of  $y^T(Xw)_+$  in  $B$ . We first consider the case where  $y^T(Xw)_+ > 0$ .

If there exists  $n \in [N]$  such that  $w; x_n \geq 0$  and  $y_n = 1$ . Consider  $v = x_n - kx_n k_2$  and let  $w = \frac{w+v}{kw+v k_2}$ , where  $\epsilon > 0$ . For index  $n \in [N]$  such that  $y_n = 1$ , as the dataset is orthogonal separable, we have  $x_n^T x_n > 0$  and

$$x_n^T(w+v) = x_n^T w + \frac{1}{kx_n k_2} x_n^T x_n > x_n^T w \quad (116)$$

This implies that  $(x_n^T w)_+ < (x_n^T(w+v))_+$ . For  $y_n = 1$ , as the data is orthogonal separable, we note that  $x_n^T x_n > 0$  and

$$x_n^T(w+v) = x_n^T w + \frac{1}{kx_n k_2} x_n^T x_n > x_n^T w \quad (117)$$

This implies that  $(x_j^T w)_+ < (x_j^T(w+v))_+$ . In summary, we have

$$y^T(X(w+v))_+ = \sum_{n=1}^N y_n (x_j^T(w+v))_+ > \sum_{n=1}^N y_n (x_j^T w)_+ = y^T(Xw)_+ > 0 \quad (118)$$

If  $w; x_n \leq 0$ , then  $w^T v < 0$ . This implies that with sufficiently small  $\epsilon$ , we have  $kw+v k_2 < kw k_2 = 1$ . Therefore,

$$y^T(Xw)_+ = \frac{1}{kw+v k_2} y^T(X(w+v))_+ > y^T(X(w+v))_+ > y^T(Xw)_+; \quad (119)$$

which leads to a contradiction. If  $w; x_n = 0$ , we note that

$$(x_n^T(w+v))_+ > (x_n^T w)_+ \quad (120)$$

This implies that

$$y^T(X(w+v))_+ > y^T(Xw)_+ \quad (121)$$

We also note that  $kw+v k_2 = \frac{1}{1+\epsilon^2} = 1 - O(\epsilon^2)$ . Therefore, with sufficiently small  $\epsilon$ , we have

$$y^T(Xw)_+ > \frac{y^T(Xw)_+}{1+\epsilon^2} > y^T(Xw)_+ \quad (122)$$

We then consider the case where  $y^T(Xw)_+ < 0$ . Apparently, we can make  $y^T(Xw)_+$  larger by replacing  $w$  by  $(1-\epsilon)w$ , where  $\epsilon \in (0, 1)$ , which leads to a contradiction.

Finally, we consider the case where  $y^T(Xw)_+ = 0$ . This implies that

$$\sum_{n: y_n=1} (x_j^T w)_+ = \sum_{n: y_n=1} (x_j^T(w+v))_+ \quad (123)$$

As  $(Xw)_+ \leq 0$ , this implies that there exists at least for one index  $n \in [N]$  such that  $y_n = 1$  and  $x_n^T w > 0$ . Let  $v = x_n - kx_n k_2$ . We note that  $\frac{1}{kw+v k_2} y^T(X(w+v))_+ > 0$  for  $\epsilon > 0$ . This leads to a contradiction.

## D.6 PROOF OF PROPOSITION 7

It is sufficient to consider the case of the local maximizer. Define  $f(\mathbf{u}; \mathbf{y}) = \frac{1}{2} \mathbf{u}^T \mathbf{X} \mathbf{y} - \frac{1}{2} \mathbf{u}^T \mathbf{X} \mathbf{X}^T \mathbf{u}$ . For  $\mathbf{u} \in \mathbb{R}^N$ , we say  $\mathbf{u}$  is a local maximizer if for all index  $n \in [N]$  with  $u_n \neq 0$ ,  $\frac{\partial f}{\partial u_n} = 0$ . We say  $\mathbf{u}$  is open if  $u_n \neq 0$  for  $n \in [N]$ . Define

$$\mathbf{S} = \mathbf{f}(\mathbf{u}; \mathbf{y}) \text{sign}(\mathbf{X}\mathbf{u}) = \mathbf{g}(\mathbf{u}; \mathbf{y}) \quad (124)$$

We start with the two lemmas.

**Lemma 8** Let  $\mathbf{u}_0 \in \mathbb{R}^N$ . Suppose that  $\mathbf{u}_0$  satisfies that  $\mathbf{u}_0 = \frac{\mathbf{g}(\mathbf{u}_0; \mathbf{y})}{\|\mathbf{g}(\mathbf{u}_0; \mathbf{y})\|_2}$ . Let  $\mathbf{u} = \text{sign}(\mathbf{u}_0)$ . Then,  $\mathbf{v} \in \mathcal{B}_2$  is a local maximizer of  $\mathbf{f}(\mathbf{X}\mathbf{u})_+$  in  $\mathcal{B}_2$  if for any open  $\mathbf{u}^0$  satisfying  $\mathbf{u}^0 = \mathbf{u}$ , we have  $\|\mathbf{g}(\mathbf{u}^0; \mathbf{y})\|_2 = \|\mathbf{g}(\mathbf{u}; \mathbf{y})\|_2$ .

**PROOF** Suppose that  $\mathcal{U}$  is open. Then  $\mathcal{S}$  is an open set. In a small neighborhood  $\mathcal{U} = \frac{\mathbf{g}(\mathbf{u}; \mathbf{y})}{\|\mathbf{g}(\mathbf{u}; \mathbf{y})\|_2} = \frac{\mathbf{g}(\mathbf{u}; \mathbf{y})}{\|\mathbf{g}(\mathbf{u}; \mathbf{y})\|_2}$ ,  $\mathbf{f}(\mathbf{X}\mathbf{u})_+ = \mathbf{u}^T \mathbf{g}(\mathbf{u}; \mathbf{y})$  is a linear function of  $\mathbf{u}$ . The Riemannian gradient of  $\mathbf{u}^T \mathbf{g}(\mathbf{u}; \mathbf{y})$  at  $\mathbf{u}$  is zero. This implies that  $\mathbf{u}$  locally maximizes  $\mathbf{f}(\mathbf{X}\mathbf{u})_+$ .

Suppose that there exists at least one zero in  $\mathbf{u}_0$ . Consider any  $\mathbf{v} \in \mathcal{B}_2$  satisfying  $\mathbf{u}_0^T \mathbf{v} = 0$ . Let  $\epsilon > 0$  be a small constant such that for any  $\mathbf{u} \in (0, \epsilon]$ ,  $\mathbf{u}_0 + \epsilon \mathbf{v} \in \mathcal{S}$  where  $\mathbf{u}_0 = \frac{\mathbf{u}_0 + \epsilon \mathbf{v}}{\|\mathbf{u}_0 + \epsilon \mathbf{v}\|_2}$ .

Suppose that  $\|\mathbf{g}(\mathbf{u}_0; \mathbf{y})\|_2 < \|\mathbf{g}(\mathbf{u}; \mathbf{y})\|_2$  for all open  $\mathbf{u}^0$  satisfying  $\mathbf{u}^0 = \mathbf{u}$ . For any  $\mathbf{u}^0$  with  $\mathbf{u}^0 = \mathbf{u}$ , we construct  $\mathbf{u}^0$  by  $u_n^0 = 1$  for  $n \in [N]$  such that  $u_n = 0$  and  $u_n^0 = u_n$  for  $n \in [N]$  such that  $u_n = 0$ . We note that  $\|\mathbf{g}(\mathbf{u}^0; \mathbf{y})\|_2 < \|\mathbf{g}(\mathbf{u}; \mathbf{y})\|_2$ . Thus,  $\|\mathbf{g}(\mathbf{u}^0; \mathbf{y})\|_2 < \|\mathbf{g}(\mathbf{u}; \mathbf{y})\|_2$ . As  $\mathbf{u}^0 \text{sign}(\mathbf{X}\mathbf{u}_s)_+ = \mathbf{u}^0 \text{sign}(\mathbf{u}_s) = \mathbf{u}_s$ , we have  $\mathbf{u}^0 \text{sign}(\mathbf{X}\mathbf{u}_s)_+ = \mathbf{u}_s = \mathbf{u} \text{sign}(\mathbf{X}\mathbf{u}_s)_+ = \mathbf{u} \text{sign}(\mathbf{u}_s) = \mathbf{u}_s$ . Therefore,  $\mathbf{u}$  is a local maximizer of  $\mathbf{f}(\mathbf{X}\mathbf{u})_+$ .

**Lemma 9** Suppose that the dataset is orthogonal separable. Let  $\mathbf{u}_0 \in \mathbb{R}^N$  satisfy that  $\text{diag}(\mathbf{y}) = \mathbf{0}$ . Suppose that  $\mathbf{u}_0$  satisfies that  $\mathbf{u}_0 = \frac{\mathbf{g}(\mathbf{u}_0; \mathbf{y})}{\|\mathbf{g}(\mathbf{u}_0; \mathbf{y})\|_2}$ . Then, for any  $\mathbf{u}^0$  satisfying  $\mathbf{u}^0 = \mathbf{u}$ , we have  $\|\mathbf{g}(\mathbf{u}^0; \mathbf{y})\|_2 = \|\mathbf{g}(\mathbf{u}; \mathbf{y})\|_2$ .

**PROOF** If there exists  $n \in [N]$  such that  $u_n = 1$  and  $y_n = -1$ , as the data is orthogonal separable, we note that

$$\mathbf{x}_n^T \mathbf{g}(\mathbf{u}; \mathbf{y}) = \mathbf{x}_n^T \mathbf{X} \mathbf{y} = \sum_{n^0: n^0 > 0} \mathbf{x}_n^T \mathbf{X}_{n^0} \mathbf{y}_{n^0} = \sum_{n^0: n^0 > 0} \mathbf{x}_n^T \mathbf{X}_{n^0} \mathbf{y}_{n^0} > 0; \quad (125)$$

which contradicts with  $\text{sign}(\mathbf{x}_n^T \mathbf{g}(\mathbf{u}; \mathbf{y})) = \text{sign}(\mathbf{x}_n^T \mathbf{u}_0) = u_n = 1$ .

Suppose that there exists  $n_1 \in [N]$  such that  $u_{n_1} = 1$  and  $y_{n_1} = 1$ . Then, as the dataset is orthogonal separable, then, for index  $n_2 \in [N]$  such that  $u_{n_2} = 0$ , we note that  $y_{n_2} < 1$ . Otherwise,

$$\mathbf{x}_{n_1}^T \mathbf{g}(\mathbf{u}; \mathbf{y}) = \mathbf{x}_{n_1}^T \mathbf{X} \mathbf{y} = \sum_{n_2: n_2 > 0} \mathbf{x}_{n_1}^T \mathbf{X}_{n_2} \mathbf{y}_{n_2} = \sum_{n_2: n_2 > 0} \mathbf{x}_{n_1}^T \mathbf{X}_{n_2} \mathbf{y}_{n_2} > 0; \quad (126)$$

which contradicts with  $\text{sign}(\mathbf{x}_{n_1}^T \mathbf{g}(\mathbf{u}; \mathbf{y})) = \text{sign}(\mathbf{x}_{n_1}^T \mathbf{u}_0) = u_{n_1} = 1$ . This also implies that the index set  $\{n \in [N] \mid y_n > 0\}$  include all data with  $u_n = 1$ .

If there exists  $n^0$  such that  $u_{n^0} = 0$  and  $\|\mathbf{g}(\mathbf{u}^0; \mathbf{y})\|_2 > \|\mathbf{g}(\mathbf{u}; \mathbf{y})\|_2$ . Then, there exists at least one index  $n \in [N]$  such that  $u_n = 0$  and  $u_n^0 = 1$ . However, from the previous derivation, we note that  $y_n = -1$  and

$$\mathbf{x}_n^T \mathbf{g}(\mathbf{u}^0; \mathbf{y}) = \mathbf{x}_n^T \mathbf{X} \mathbf{y} = \sum_{j: u_j^0 > 0} \mathbf{x}_n^T \mathbf{X}_j \mathbf{y}_j = \sum_{n_1: u_{n_1}^0 > 0} \mathbf{x}_n^T \mathbf{X}_{n_1} \mathbf{y}_{n_1} < 0; \quad (127)$$

which contradicts with  $u_n^0 = 1$ .

By combining Lemma 8 and 9, we complete the proof.

## D.7 PROOF OF THEOREM 4

PROOF For almost all initialization, we can find two neurons such that  $\text{sign}(w_{2;i_+}) = \text{sign}(y^T(Xw_{1;i_+})_+) = 1$  and  $\text{sign}(w_{2;i_-}) = \text{sign}(y^T(Xw_{1;i_-})_+) = -1$  at initialization. By choosing a sufficiently small  $\epsilon > 0$  in Proposition 6, there exist two neurons  $w_{1;i_+}, w_{1;i_-}$  and times  $T_+, T_- > 0$  such that  $\cos(w_{1;i_+}(T_+); g(w_{1;i_+}(T_+); y)) > 1 - \epsilon$  and  $\cos(w_{1;i_-}(T_-); g(w_{1;i_-}(T_-); y)) < -(1 - \epsilon)$ . This implies that  $w_{1;i_+}(T_+)$  and  $w_{1;i_-}(T_-)$  are sufficiently close to certain stationary points of gradient flow maximizing/minimizing  $\langle w, Xu_+ \rangle$  over  $B$ , i.e.,  $f(u) = \langle u, g(u; y) \rangle = 1$ . As the dataset is orthogonal separable, according to Lemma 4 and Proposition 7, the corresponding diagonal matrices  $D_{i_+}(T_+)$  and  $D_{i_-}(T_-)$  satisfy that  $D_{i_+}(T_+) = \text{diag}(I(y = 1))$  and  $D_{i_-}(T_-) = \text{diag}(I(y = -1))$ . According to Lemma 3 in (Phuong & Lampert, 2021), we have  $\hat{D}_{i_+}(t) = \text{diag}(I(y = 1))$  and  $\hat{D}_{i_-}(t) = \text{diag}(I(y = -1))$  hold for  $t \geq \max\{T_+, T_-\}$ .

With  $t \geq 1$ , according to Proposition 4, the dual variable  $\lambda$  in the KKT point of the non-convex max-margin problem (13) is dual feasible, i.e., satisfies (16). Suppose that  $\lambda$  is a limiting point of  $\frac{\lambda(t)}{\kappa \frac{(t)}{(t)k_2} t_0}$  and  $\mu$  is the corresponding dual variable. From Theorem 1, we note that the pair  $(\lambda; \mu)$  corresponds to the KKT point of the convex max-margin problem (14).

## E PROOFS OF MAIN RESULTS ON MULTICLASS CLASSIFICATION

### E.1 PROOF OF PROPOSITION 1

The neural network training problem (4) can be separated into  $k$  subproblems. Each of these subproblems corresponds to the neural network training problem (3) for binary classification. For each subproblem, by applying Proposition 2, we complete the proof.

### E.2 PROOF OF THEOREM 1

We note that the neural network training problem (4) can be separated into  $k$  subproblems. Each of these subproblems corresponds to the neural network training problem (3) for binary classification. By applying Proposition 6 with  $\mu = y_k$  to each subproblem with  $\mu = y_k$ , we complete the proof.

### E.3 PROOF OF THEOREM 2

Similarly, the corresponding non-convex max-margin problem (5) and the convex max-margin problem (7) can be separated into  $k$  subproblems. Each of these subproblems corresponds to the non-convex max-margin problem (2) and the convex max-margin problem (1) for binary classification. By applying Theorem 4 to each subproblem with  $\mu = y_k$ , we complete the proof.

## F NUMERICAL EXPERIMENT

### F.1 DETAILS ON FIGURE 5

We provide the experiment setting in Figure 1 and 5 as follows. The dataset is given by  $\begin{matrix} 1:65 & 0:47 \\ 0:47 & 1:35 \end{matrix} \in \mathbb{R}^2 \times \mathbb{R}^2$  and  $y = \begin{matrix} 1 \\ 1 \end{matrix} \in \mathbb{R}^2$ . Here we have  $N = 2$  and  $d = 2$ . We note that this dataset is orthogonal separable but not spike-free. We plot the ellipsoid set and the rectified ellipsoid set in Figure 6.



Figure 6: The ellipsoid set and the rectified ellipsoid set. Orthogonal separable dataset.

We enumerate all possible hyperplane arrangements in the and solve the convex max-margin problem (14) via CVXPY to obtain the following non-zero neurons

$$u_{1;3} = \begin{pmatrix} 0.58 \\ 0.16 \end{pmatrix}; \quad w_{1;2}^0 = \begin{pmatrix} 0.23 \\ 0.66 \end{pmatrix}; \quad (128)$$

We note that the dual problem (15) is equivalent to

$$\begin{aligned} \max \quad & \sum_j y_j; \\ \text{s.t.} \quad & k \|X^T D_j - X^T (2D_j - I) z_{j,+}\|_2 \leq 1; \quad j = 1, \dots, 2 [p]; \\ & k \|X^T D_j - X^T (2D_j - I) z_{j,-}\|_2 \leq 1; \quad j = 1, \dots, 2 [p]; \\ & z_{j,+} \geq 0; \quad z_{j,-} \leq 0; \quad j = 1, \dots, 2 [p]; \quad \text{diag}(y) \geq 0. \end{aligned} \quad (129)$$

The above problem is a second-order cone program (SOCP) and can be solved via standard convex optimization frameworks such as CVX and CVXPY. We solve (129) to obtain the optimal dual variable  $\lambda$ . For the geometry of the dual problem, as the dataset is orthogonal separable, the set  $f^* = \max_{\|u\|_2 \leq 1} \sum_j y_j \langle X u, u \rangle$  reduces to  $f^* = \max_{\|u\|_2 \leq 1} \sum_j y_j \langle X u_1, u_1 \rangle + \sum_j y_j \langle X u_2, u_2 \rangle$ , where  $u_1, u_2$  correspond to two vectors at the spikes of the rectified ellipsoid set. We draw the sets  $f^*$ , the optimal dual variable and the direction of  $f^*$  in Figure 2.

For each  $D_j \in \mathcal{P}$ , we solve for the vector  $u_j$  which maximize/minimize  $D_j X u_j$  with the constraints  $\|u_j\|_2 \leq 1$  and  $(2D_j - I) X u_j = 0$ . We plot the rectified ellipsoid set  $\{X u_j\}_j$ , vectors  $u_j$ , neurons in the optimal solution (14) scaled to unit  $\ell_2$ -norm and the direction of  $f^*$  in Figure 1. We note that each neuron in the optimal solution from (14) (scaled to unit  $\ell_2$ -norm) maximize/minimize the corresponding  $D_j X u_j$  given  $(2D_j - I) X u_j = 0$ .

Then, we consider a two-layer ReLU network with  $m = 10$  neurons and apply the gradient descent method to train on the logistic loss (3). Let  $w_{1;i} = \frac{w_{1;i}}{\|w_{1;i}\|_2}$  for  $i = 1, \dots, m$ . We plot  $w_{1;i}$  and  $\langle X w_{1;i}, u_j \rangle$  at iteration  $t = 0, \dots, 4g$  along with neurons in the optimal solution (14) scaled to unit  $\ell_2$ -norm in Figure 5. Certain neurons do not move, while the activated neurons trained by gradient descent tend to converge to the direction of the neurons in the optimal solution to (14).

We repeat the training on the logistic loss (3) with the gradient descent method several times and we plot the trajectories in Figure 7.

Figure 7: Multiple independent random initializations of gradient descent trajectories on the same orthogonal separable dataset.

## F.2 EXPERIMENT ON SPIKE-FREE DATASET

We repeat the previous numerical experiment on a non-spike-free dataset:  $\begin{matrix} 1:65 & 0:47 \\ 0:47 & 1:35 \end{matrix}$   $\mathbb{R}^2$  and  $\mathbb{R}^2$ . Similarly, we plot the ellipsoid set and the rectified set in Figure 8.

Figure 8: The ellipsoid set and the rectified ellipsoid set for a non-spike-free dataset.

We enumerate all possible hyperplane arrangements in  $\mathbb{R}^2$  and solve the convex max-margin problem (14) via CVXPY to obtain the following non-zero neuron

$$u_{1:4} = \begin{matrix} 0:43 \\ 0:59 \end{matrix} \quad (130)$$

We plot the rectified ellipsoid set  $(Xu)_+$   $\|u\|_2 = 1$ , vectors  $u_j$ , neurons in the optimal solution to (14) scaled to unit  $\ell_2$ -norm and the direction of  $u_{1:j}$  and  $(Xu_{1:j})_+$  at iteration  $10^j$   $j = 0, \dots, 4$  along with neurons in the optimal solution (14) scaled to unit  $\ell_2$ -norm in Figure 10.

Figure 9: Rectified Ellipsoidal set and corresponding extreme points for a non-spike-free dataset.

