

When Are Bias-Free ReLU Networks Effectively Linear Networks?

Anonymous authors

Paper under double-blind review

Abstract

We investigate the implications of removing bias in ReLU networks regarding their expressivity and learning dynamics. We first show that two-layer bias-free ReLU networks have limited expressivity: the only odd function two-layer bias-free ReLU networks can express is a linear one. We then show that, under symmetry conditions on the data, these networks have the same learning dynamics as linear networks. This enables us to give analytical time-course solutions to certain two-layer bias-free (leaky) ReLU networks outside the lazy learning regime. While deep bias-free ReLU networks are more expressive than their two-layer counterparts, they still share a number of similarities with deep linear networks. These similarities enable us to leverage insights from linear networks to understand certain ReLU networks. Overall, our results show that some properties previously established for bias-free ReLU networks arise due to equivalence to linear networks.

1 Introduction

Theorists make simplifications to real-world models because simplified models are mathematically more tractable, yet discoveries made in them may hold in general. For instance, linear models have illuminated benign overfitting (Bartlett et al., 2020) and double descent (Advani et al., 2020) in practical neural networks. In this paradigm, understanding the consequences of a simplification is critical, since it informs us which discoveries in simple models extend to complex ones. Here we inspect a specific simplification that is common in theoretical work on ReLU networks (Zhang et al., 2019; Du et al., 2019; Arora et al., 2019; Lyu & Li, 2020; Vardi & Shamir, 2021): the removal of the bias terms. The removal of bias not only appears in theoretical work but also has practical applications. Some real-world models adopt bias removal to introduce scale invariance, which can be a beneficial property for image denoising (Mohan et al., 2020; Zhang et al., 2022), image classification (Zarka et al., 2021), and diffusion models (Kadkhodaie et al., 2024). This paper seeks to illuminate the implications of bias removal in ReLU networks, and so provide insight for theorists on when bias removal is desirable.

We investigate how removing bias affects the expressivity and the learning dynamics of ReLU networks and identify scenarios where bias-free ReLU networks are effectively linear networks. For expressivity, we show that two-layer bias-free (leaky) ReLU networks cannot express odd functions except linear functions. This was proven for input uniformly distributed on a sphere (Basri et al., 2019, Theorem 2 and 4), but we prove it for arbitrary input with a simpler approach. We then consider deep bias-free (leaky) ReLU networks and show a depth separation result, i.e., deep bias-free ReLU networks can express homogeneous nonlinear odd functions while two-layer ones cannot. For learning dynamics, we show that two-layer bias-free (leaky) ReLU networks have the same learning dynamics as a linear network when trained with square loss or logistic loss on symmetric datasets, whose target function is odd and input distribution is even. Our symmetry Condition 3 on the dataset incorporates the datasets studied in several prior works (Sarussi et al., 2021; Lyu et al., 2021; Zhang et al., 2024). We also present two cases where two-layer bias-free ReLU networks evolve like multiple independent linear networks. Finally, we empirically find that when the target function is linear, deep bias-free ReLU networks form low-rank weights similar to those in deep linear networks.

By revealing regimes where bias-free ReLU networks behave like linear networks, we provide an accessible way of understanding ReLU networks within these regimes, as well as a cautionary note that studying nonlinear behaviors generally requires stepping beyond these regimes. This perspective leverages insights from linear networks, which are simpler and thus enjoy much richer theoretical results than ReLU networks (Baldi & Hornik, 1989; Fukumizu, 1998; Saxe et al., 2014; 2019; Arora et al., 2018; Ji & Telgarsky, 2019; Lampinen & Ganguli, 2019; Gidel et al., 2019; Tarmoun et al., 2021; Braun et al., 2022; Ziyin et al., 2022). For example, we are able to give closed-form time-course solutions to certain two-layer ReLU networks outside the lazy learning regime in Corollary 8. Our findings suggest that the bias terms in a ReLU network play an important role in learning nonlinear tasks. Our contributions are the following:

- Section 3 proves the limited expressivity of bias-free (leaky) ReLU networks, and shows a depth separation result between two-layer and deep bias-free ReLU networks;
- Section 4.1 proves that, under symmetry Condition 3 on the dataset, two-layer bias-free (leaky) ReLU networks trained with square loss or logistic loss evolve the same as linear networks, and gives analytical time-course solutions for ReLU networks in this regime.
- Section 4.2 shows that bias-free ReLU networks behave similarly to multiple independent linear networks on orthogonal and XOR datasets;
- Section 5 shows the similarities between deep bias-free ReLU networks and deep linear networks, and finds specific rank-one and rank-two structure in the weights.

1.1 Related Work

Basri et al. (2019) proved, using harmonic analysis, that two-layer bias-free ReLU networks can neither learn nor express odd nonlinear functions when input is uniformly distributed on a sphere (Basri et al., 2019, Theorem 2 and 4). We make a similar argument with a simpler proof. Our Theorem 1 handles arbitrary input, includes both ReLU and leaky ReLU networks, and the proof only involves rewriting the (leaky) ReLU activation function as the sum of a linear function and an absolute value function.

Lyu et al. (2021) proved two-layer bias-free leaky ReLU networks trained with logistic loss converge to a linear, max-margin classifier on linearly separable tasks with a data augmentation procedure. Our Theorem 7 shows that the learning dynamics of leaky ReLU networks in their setup is equivalent to that of a linear network. In light of this equivalence, their result is guaranteed given that linear networks trained with logistic loss converge to the max-margin classifier on linearly separable tasks (Soudry et al., 2018). In addition, we relax the assumption on the task from being linearly separable to being odd, and thus identify a practical challenge: the data augmentation procedure of Lyu et al. (2021) can cause the ReLU network to fail to learn a linearly non-separable task — a task the network might have succeeded to learn without data augmentation.

Zhang et al. (2024) found that two-layer bias-free ReLU networks have similar loss and weight norm curves as linear networks when trained on datasets with zero mean Gaussian input and a linear target. They reported that training the ReLU networks is about twice as slow as their linear counterpart. Our Theorem 7 explains their observation: we prove that the dynamics of two-layer bias-free ReLU networks is exactly twice as slow as their linear counterpart for a general class of datasets, including theirs.

A few other works have alluded to the connections between two-layer ReLU and linear networks. Sarussi et al. (2021) discovered that two-layer bias-free leaky ReLU networks converge to a decision boundary that is very close to linear when the teacher model is linear. Their theoretical results assume that the second layer is fixed while we train all layers of the network. Saxe et al. (2022) studied gated deep linear networks and found they closely approximate a two-layer bias-free ReLU network trained on an XOR task. But they did not generalize the connection between gated linear networks and ReLU networks beyond the XOR case. Boursier & Flammarion (2024) gave an example dataset with three scalar input data points, in which two-layer bias-free (leaky) ReLU networks converge to the linear, ordinary least square estimator. Holzmüller & Steinwart (2022) studied two-layer leaky ReLU networks with bias and found that they perform linear regression on certain data distributions with scalar input, because the bias fails to move far away from their initialization at zero.

While prior works have studied cases where bias-free ReLU networks behave like linear networks, their connections have not been explicitly highlighted or systematically summarized. This paper aims to bring the connections between bias-free ReLU and linear networks into focus, offering a comparative perspective on ReLU networks.

2 Preliminaries

Notation: We use bold symbols to denote vectors and matrices. Double-pipe brackets $\|\cdot\|$ denote the L2 norm of a vector or the Frobenius norm of a matrix. Angle brackets $\langle\cdot\rangle$ denote the average over the dataset. The circled dot \odot denotes the element-wise product.

2.1 Two-Layer Bias-Free (Leaky) ReLU and Linear Networks

A two-layer bias-free (Leaky) ReLU network with H hidden neurons is defined as

$$f(\mathbf{x}; \mathbf{W}) = \mathbf{W}_2 \sigma(\mathbf{W}_1 \mathbf{x}) = \sum_{h=1}^H w_{2h} \sigma(w_{1h} \mathbf{x}), \quad \text{where } \sigma(z) = \max(z, \alpha z), \alpha \in [0, 1]. \quad (1)$$

Here $\mathbf{x} \in \mathbb{R}^D$ is the input, $\mathbf{W}_1 \in \mathbb{R}^{H \times D}$ is the first-layer weight, $\mathbf{W}_2 \in \mathbb{R}^{1 \times H}$ is the second-layer weight, and \mathbf{W} denotes all weights collectively. This is a ReLU network when $\alpha = 0$ and a leaky ReLU network when $\alpha \in (0, 1)$. When $\alpha = 1$, the network is a linear network, and can be written as $f(\mathbf{x}) = \mathbf{W}_2 \mathbf{W}_1 \mathbf{x}$. We also denote the linear network as $f^{\text{lin}}(\mathbf{x}; \mathbf{W}^{\text{lin}}) = \mathbf{W}_2^{\text{lin}} \mathbf{W}_1^{\text{lin}} \mathbf{x}$ when we need to distinguish it from ReLU networks.

We consider the rich regime (Woodworth et al., 2020) in which the network is initialized with small random weights. The network is trained with gradient descent on a dataset $\{\mathbf{x}_\mu, y_\mu\}_{\mu=1}^P$ consisting of P samples. We study square loss $\mathcal{L} = \langle (y - f(\mathbf{x}))^2 \rangle / 2$ and logistic loss $\mathcal{L}_{\text{LG}} = \langle \ln(1 + e^{yf(\mathbf{x})}) \rangle$. We focus on square loss in the main text and provide derivations with logistic loss in the appendix. The learning rate is η and the inverse of the learning rate is the time constant $\tau = 1/\eta$. In the limit of small learning rate, the gradient descent dynamics are well approximated by the gradient flow differential equations

$$\tau \dot{\mathbf{W}}_1 = \langle \sigma'(\mathbf{W}_1 \mathbf{x}) \odot \mathbf{W}_2^\top (y - \mathbf{W}_2 \sigma(\mathbf{W}_1 \mathbf{x})) \mathbf{x}^\top \rangle, \quad (2a)$$

$$\tau \dot{\mathbf{W}}_2 = \langle (y - \mathbf{W}_2 \sigma(\mathbf{W}_1 \mathbf{x})) \sigma(\mathbf{W}_1 \mathbf{x})^\top \rangle, \quad (2b)$$

where σ' is the derivative of σ , \odot is the element-wise product, and the angle brackets $\langle\cdot\rangle$ denote taking the average over the dataset.

For linear networks, $\sigma(z) = z$, the gradient flow dynamics can be written as

$$\tau \dot{\mathbf{W}}_1^{\text{lin}} = \mathbf{W}_2^{\text{lin}^\top} (\beta^\top - \mathbf{W}_2^{\text{lin}} \mathbf{W}_1^{\text{lin}} \Sigma), \quad (3a)$$

$$\tau \dot{\mathbf{W}}_2^{\text{lin}} = (\beta^\top - \mathbf{W}_2^{\text{lin}} \mathbf{W}_1^{\text{lin}} \Sigma) \mathbf{W}_1^{\text{lin}^\top}, \quad (3b)$$

where Σ denotes the input data covariance and β denotes the input-output correlation,

$$\Sigma = \langle \mathbf{x} \mathbf{x}^\top \rangle, \quad \beta = \langle y \mathbf{x} \rangle. \quad (4)$$

2.2 Deep Networks

A deep neural network of depth L is $f(\mathbf{x}) = h_L$ where h_L is recursively defined as

$$\begin{aligned} \mathbf{h}_l &= \mathbf{W}_l \sigma(\mathbf{h}_{l-1}), \quad 2 \leq l \leq L, \\ \mathbf{h}_1 &= \mathbf{W}_1 \mathbf{x}. \end{aligned} \quad (5)$$

Here $\mathbf{h}_1, \dots, \mathbf{h}_{L-1}$ are vectors and h_L is the scalar output. The gradient flow dynamics trained with square loss is

$$\tau \dot{\mathbf{W}}_l = \left\langle \frac{\partial h_L}{\partial \mathbf{h}_l} (y - h_L) \sigma(\mathbf{h}_{l-1})^\top \right\rangle. \quad (6)$$

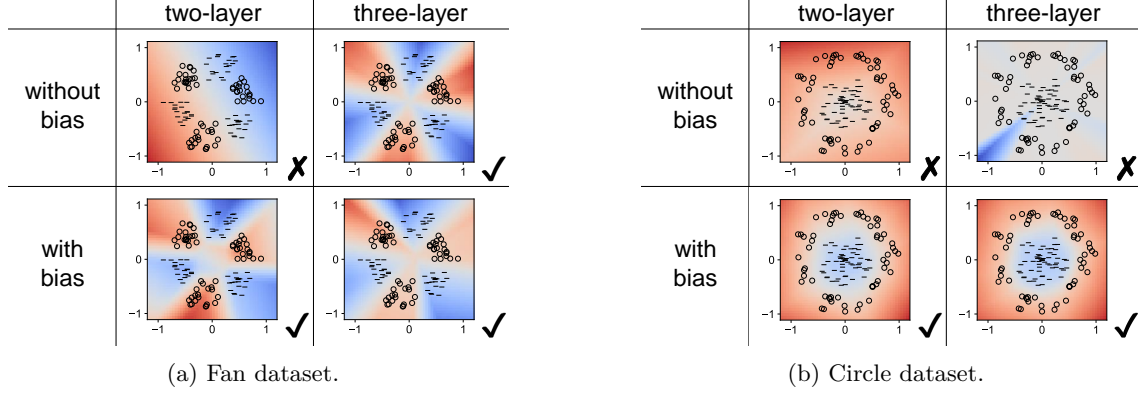


Figure 1: The expressivity of two-layer and deep ReLU networks with and without bias. The networks are trained with logistic loss until the loss stops decreasing. The empty circles are data points with $+1$ labels; short lines are data points with -1 labels. The network output is plotted in color. (a) The fan dataset is odd, homogeneous, and satisfies Condition 3. Two-layer bias-free ReLU networks cannot express it. (b) The circle dataset is not homogeneous. Two-layer and deep bias-free ReLU networks cannot express it. Experimental details are provided in Appendix H.

For deep linear networks, the gradient flow dynamics can be written as

$$\tau \dot{\mathbf{W}}_l^{\text{lin}} = \left(\prod_{i=l+1}^L \mathbf{W}_i^{\text{lin}} \right)^\top \left(\boldsymbol{\beta}^\top - \prod_{i=1}^L \mathbf{W}_i^{\text{lin}} \boldsymbol{\Sigma} \right) \left(\prod_{i=1}^{l-1} \mathbf{W}_i^{\text{lin}} \right)^\top, \quad (7)$$

where $\prod_i \mathbf{W}_i$ represents the ordered product of matrices with the largest index on the left and smallest on the right.

3 Network Expressivity

We first examine the expressivity of bias-free ReLU networks. It is well known that standard ReLU networks are universal approximators (Hornik et al., 1989; Pinkus, 1999) while bias-free ReLU networks are not since they can only express positively homogeneous functions, i.e., $g(a\mathbf{x}) = ag(\mathbf{x}) \forall a > 0$. Moreover, Section 3.1 shows that two-layer bias-free ReLU networks cannot express any odd function except linear functions. Section 3.2 shows that deep bias-free ReLU networks are more expressive than two-layer ones, but are still limited to positively homogeneous functions.

3.1 Two-Layer Bias-Free (Leaky) ReLU Networks

Theorem 1. *Two-layer bias-free (leaky) ReLU networks can only express a linear function plus a positively homogeneous even function.*

Proof. An arbitrary two-layer (leaky) ReLU network can be written as

$$\sum_{h=1}^H w_{2h} \sigma(\mathbf{w}_{1h} \mathbf{x}) = \sum_{h=1}^H w_{2h} \left[\frac{1+\alpha}{2} \mathbf{w}_{1h} \mathbf{x} + \frac{1-\alpha}{2} |\mathbf{w}_{1h} \mathbf{x}| \right], \quad (8)$$

$$\text{---} \diagup = \text{---} \diagup + \text{---} \diagdown$$

which is a linear function plus a positively homogeneous even function. \square

Corollary 2. *The only odd function that bias-free two-layer (leaky) ReLU networks can express is the linear function.*

Due to this restricted expressivity, two-layer bias-free ReLU networks fail to classify the fan dataset, as shown in Figure 1a.

3.2 Deep Bias-Free (Leaky) ReLU Networks

Similarly to two-layer bias-free ReLU networks, deep bias-free ReLU networks can express only positively homogeneous functions. Thus, as shown in Figure 1b, both two-layer and deep bias-free ReLU networks fail to classify the circle dataset. However, in contrast to two-layer bias-free ReLU networks, deep bias-free ReLU networks can express some odd nonlinear functions. For instance, for two-dimensional input $\mathbf{x} = [x_1, x_2]^\top$, the function below is odd, nonlinear, and can be implemented by a three-layer bias-free ReLU network,

$$g(\mathbf{x}) = \sigma(\sigma(x_1) - \sigma(x_2)) - \sigma(\sigma(-x_1) - \sigma(-x_2)), \quad \text{where } \sigma(z) = \max(z, 0). \quad (9)$$

Thus, we have a depth separation result for bias-free ReLU networks: there exist odd nonlinear functions, such as $g(\mathbf{x})$ defined in Equation (9) and visualized in Figure 8, that two-layer bias-free ReLU networks cannot express but deep bias-free ReLU networks can.

4 Learning Dynamics in Two-Layer Bias-Free ReLU Networks

4.1 Symmetric Datasets

Section 3.1 has proven that two-layer bias-free ReLU networks cannot express odd functions except linear functions. We now show that under the Condition 3 on the dataset, two-layer bias-free ReLU networks not only find a linear solution but also have the same learning dynamics as a two-layer linear network.

Condition 3. The dataset satisfies the following two symmetry conditions:

1. The empirical input data distribution is even: $p(\mathbf{x}) = p(-\mathbf{x})$;
2. The target model is odd: $y(\mathbf{x}) = -y(-\mathbf{x})$.

Remark 4. For infinite data, the first part of Condition 3 includes common distributions, such as any Gaussian distribution with zero mean. For finite data, the first part of Condition 3 means that if \mathbf{x} is present in the dataset, $-\mathbf{x}$ is also present. Condition 3 includes the dataset studied in Lyu et al. (2021). They considered linearly separable binary classification tasks with a data augmentation procedure in which $(-\mathbf{x}, -y)$ is added to the dataset if (\mathbf{x}, y) is in the dataset. We have the same assumption on the input data distribution but relax the assumption on the target model from being linearly separable to being odd.

Assumption 5. At initialization, there exist an unit vector \mathbf{r} such that $\mathbf{W}_1 = \mathbf{W}_2^\top \mathbf{r}^\top$ and the second-layer weight \mathbf{W}_2 has equal L2 norms for its positive and negative elements.

Remark 6. Under Condition 3, the dynamics of two-layer bias-free (leaky) ReLU networks and linear networks initialized with small weights can both be approximated by a linear differential equation in the early phase of learning. Thus, the weights of the ReLU network and the weights of the linear network form the same rank-one structure and the left singular vector of \mathbf{W}_1 is aligned with \mathbf{W}_2 , as proven in Theorem 15. At the end of the early phase, the weights are rank-one and aligned with bounded errors. To simplify the analysis, Assumption 5 assumes the weights are exactly rank-one and aligned, which is justified by the initialization approaching 0. We also assume, for simplicity, that the second-layer weights have equal L2 norm for their positive and negative elements, which is justified by the width approaching infinity. Simulations in Figure 2b show that errors are small in the finite case as well.

Under Assumption 5 and Condition 3, the learning dynamics of two-layer bias-free (leaky) ReLU networks reduces to (see Appendix C.2)

$$\tau \dot{\mathbf{W}}_1 = \frac{\alpha + 1}{2} \mathbf{W}_2^\top \beta^\top - \left(\frac{\alpha + 1}{2} \right)^2 \mathbf{W}_2^\top \mathbf{W}_2 \mathbf{W}_1 \Sigma, \quad (10a)$$

$$\tau \dot{\mathbf{W}}_2 = \frac{\alpha + 1}{2} \beta^\top \mathbf{W}_1^\top - \left(\frac{\alpha + 1}{2} \right)^2 \mathbf{W}_2 \mathbf{W}_1 \Sigma \mathbf{W}_1^\top. \quad (10b)$$

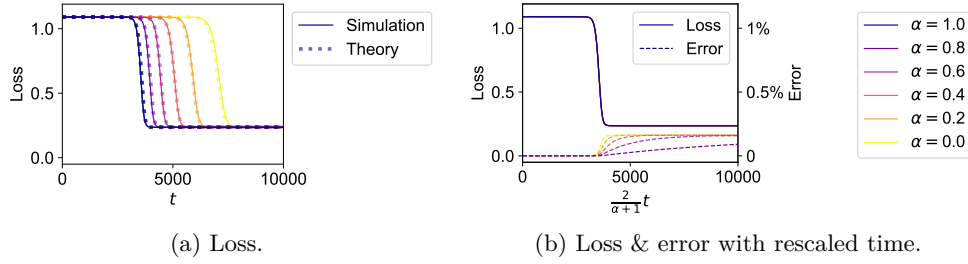


Figure 2: Two-layer bias-free (leaky) ReLU networks can evolve like a linear network. (a) Loss curves with different leaky ReLU parameter α (note $\alpha = 1$ is a linear network). The simulations match the theoretical solutions in Equation (13). The loss converges to global minimum, which is not zero due to the restricted expressivity of two-layer bias-free ReLU networks. (b) The simulated loss curves are plotted against a rescaled time axis; they collapse to one curve, demonstrating the networks are implementing the same linear function as in Equation (11). The error, defined as $\left\| \sqrt{\frac{\alpha+1}{2}} \mathbf{W} \left(\frac{2}{\alpha+1} t \right) - \mathbf{W}^{\text{lin}}(t) \right\| / \left\| \mathbf{W}^{\text{lin}}(t) \right\|$, is less than 0.3%, demonstrating that the weights in the (leaky) ReLU network are close to the weights in the linear network as in Equation (12). The errors are not exactly zero because the initialization is small but nonzero in the simulations. Experimental details are provided in Appendix H.

Except for the constant coefficients, these equations are the same as the ones for the linear network given in Equation (3). Thus, apart from the fact that learning is $(\alpha + 1)/2$ times slower and the weights are $\sqrt{2/(\alpha + 1)}$ times larger, the ReLU network has the same learning dynamics as its linear counterpart. We formally state this equivalence below.

Theorem 7. *A two-layer (leaky) ReLU network and a linear network are trained with square or logistic loss starting from weights which differ by a scale factor, $\mathbf{W}(0) = \sqrt{2/(\alpha + 1)} \mathbf{W}^{\text{lin}}(0)$. Under Condition 3 on the dataset and Assumption 5 on the initial weights, we have that $\forall t \geq 0$, Assumption 5 remains valid and:*

1. *The (leaky) ReLU network implements the same linear function as the linear network with scaled time*

$$f(\mathbf{x}; \mathbf{W}(t)) = f^{\text{lin}} \left(\mathbf{x}; \mathbf{W}^{\text{lin}} \left(\frac{\alpha + 1}{2} t \right) \right); \quad (11)$$

2. *The weights in the (leaky) ReLU network are the same as scaled weights in the linear network*

$$\mathbf{W}(t) = \sqrt{\frac{2}{\alpha + 1}} \mathbf{W}^{\text{lin}} \left(\frac{\alpha + 1}{2} t \right). \quad (12)$$

As a sanity check, setting $\alpha = 1$ yields a linear network, where \mathbf{W} is identical to \mathbf{W}^{lin} . As an example, setting $\alpha = 0$ yields a ReLU network, where learning slows down by half, explaining a previous empirical observation (Zhang et al., 2024). We validate Theorem 7 and the plausibility of Assumption 5 with numerical simulations in Figure 2. Additionally, we provide theoretical proof that Theorem 7 holds with L2 regularization and empirical evidence that some of Theorem 7 hold with large initialization and a moderately large learning rate in Appendices C.4 to C.6.

If the input covariance is white, we can further write down the exact time-course solution in closed form for two-layer bias-free (leaky) ReLU networks by adopting the solutions from linear networks (Braun et al., 2022, Theorem 3.1). This gives us the following corollary.

Corollary 8. *For learning with square loss, if the input covariance is white, $\Sigma = \mathbf{I}$, the solution to Equation (11) is $f(\mathbf{x}; \mathbf{W}) = \mathbf{w}(t)^\top \mathbf{x}$ with*

$$\mathbf{w}(t) = \left(1 + \frac{q_1}{q_2} e^{-2s\bar{t}} \right) \left[\bar{\beta} \left(1 - \frac{q_1}{q_2} e^{-2s\bar{t}} \right) + \frac{2}{q_2} (\mathbf{I} - \bar{\beta} \bar{\beta}^\top) \mathbf{r} e^{-s\bar{t}} \right] \left[\frac{4}{q_2^2} \left(w_{\text{init}}^{-2} + \left(1 - (\mathbf{r}^\top \bar{\beta})^2 \right) \bar{t} \right) e^{-2s\bar{t}} + \frac{1}{s} \left(1 + \frac{q_1^2}{q_2^2} e^{-2s\bar{t}} \right) (1 - e^{-2s\bar{t}}) \right]^{-1}, \quad (13)$$

where \tilde{t} is a shorthand for rescaled time $\tilde{t} = \frac{\alpha+1}{2\tau}t$ and the constant quantities are $s = \|\beta\|$, $\bar{\beta} = \beta/s$, $q_1 = 1 - \mathbf{r}^\top \bar{\beta}$, $q_2 = 1 + \mathbf{r}^\top \bar{\beta}$, $w_{\text{init}} = \|\mathbf{W}_1(0)\|$.

The solution given in Equation (13) matches simulations, as shown in Figure 2a.

Since the time evolution of two-layer bias-free (leaky) ReLU networks is the same as that of linear networks (modulo scale factors), their converged weights will also be the same. For learning with square loss, linear networks converge to the ordinary least squares solution (Saxe et al., 2014). For linearly separable binary classification with logistic loss, linear networks converge to the max-margin (hard margin SVM) solution (Soudry et al., 2018). Thus two-layer bias-free (leaky) ReLU networks also converge to these solutions when they behave like linear networks; see Appendix C.3.

Corollary 9. *Under the same conditions as Theorem 7, the two-layer bias-free (leaky) ReLU network converges to a linear solution $f(\mathbf{x}; \mathbf{W}(\infty)) = \mathbf{w}^*{}^\top \mathbf{x}$. For square loss, \mathbf{w}^* is the ordinary least squares solution, $\Sigma^{-1}\beta$, which is the global minimum. For linearly separable binary classification with logistic loss, \mathbf{w}^* is the max-margin solution.*

4.2 Orthogonal and XOR Datasets

In Section 4.1 we considered symmetric datasets where a two-layer bias-free ReLU network evolves like one linear network. We now present two datasets where a two-layer bias-free ReLU network evolves like multiple independent linear networks.

The first dataset has orthogonal inputs, which is a common setting studied by prior literature (Boursier et al., 2022; Telgarsky, 2023; Frei et al., 2023b;c; Kou et al., 2023). In particular, we use the same dataset as Boursier et al. (2022, Figure 3), which consists of two orthogonal data points, as shown in Figure 3a. We train a two-layer bias-free ReLU network on this dataset to reproduce the loss curve in (Boursier et al., 2022, Figure 3). We then train two two-layer linear networks on each data point separately. We find that the timing and the amount of the loss drop overlap with the loss curves of the two linear networks as shown in Figure 3b. To understand this overlap, we plot the first-layer weights of the ReLU network in black arrows in Figure 3a and find that the weights align with either one of the two data points. Since the two directions are orthogonal, the learning dynamics of the two groups of neurons decouple, as derived in Appendix D. Each group of neurons evolves like a linear network trained on that single data point. Hence, weights in the ReLU network evolve like a linear network trained on either one of the two data points separately. The same applies to learning with logistic loss, as shown in Figure 7, which was not covered in Boursier et al. (2022).

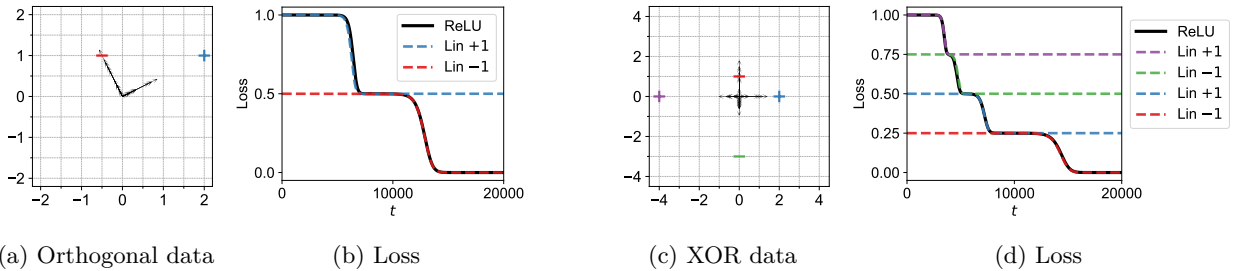


Figure 3: Two-layer bias-free ReLU networks can evolve like several linear networks. (a) An orthogonal input dataset used in (Boursier et al., 2022, Figure 3). The + and - signs represent data points with +1 and -1 labels respectively. Their different colors are used only to distinguish the loss curves. The black arrows are the first-layer weights at convergence. (b) The loss curve of the ReLU network overlaps with two linear networks trained on each of the two data points respectively. (c) An XOR-like dataset. (d) The loss curve of the ReLU network overlaps with four linear networks trained on each of the four data points separately. Details: We use summed (instead of averaged) square loss for this figure. The initial losses are vertically aligned to help illustrate the overlap. More details are in Appendix H.

We observe similar behavior in the XOR-like task shown in Figure 3c, which is neither linearly separable nor odd. XOR-like datasets are also a common setting in prior literature (Saxe et al., 2022; Frei et al., 2023a; Meng et al., 2024; Xu et al., 2024; Glasgow, 2024). As shown in Figure 3d, we find that the loss curves of a two-layer bias-free ReLU network trained on the XOR task overlap with four linear networks trained on each data point separately. In this case, the dynamics of multiple linear networks well approximate that of a ReLU network, even though the ReLU network learns a nonlinear function.

In Figures 3b and 3d, the loss curves go through multiple decreases, each corresponding to learning a data point. Similar behaviors were examined by Boursier et al. (2022); Xu et al. (2024) and characterized as saddle-to-saddle dynamics. The connections we find between ReLU and linear networks may help understand these behaviors in ReLU networks because saddle-to-saddle dynamics has been well studied for linear networks (Saxe et al., 2014; 2019; Jacot et al., 2021; Pesme & Flammarion, 2023).

5 Learning Dynamics in Deep Bias-Free ReLU Networks

In Section 4 we showed that two-layer bias-free ReLU networks behave like linear networks under symmetry Condition 3 and small initialization. This does not extend to deep bias-free networks. When trained on a dataset satisfying Condition 3, deep bias-free ReLU networks can learn nonlinear solutions if the target function is nonlinear, as shown in Figure 1a (upper right). Nonetheless, we find deep bias-free ReLU networks can form low-rank weights that are similar to those in deep linear networks. We give an example where the empirical input distribution is even and the target function is linear.

In a deep linear network, weights form an approximately rank-one structure and adjacent layers are approximately aligned when trained from small initialization (Ji & Telgarsky, 2019; Advani et al., 2020; Atanasov et al., 2022; Marion & Chizat, 2024). The rank-one weight matrices can be written approximately as outer-products of two vectors

$$\mathbf{W}_1^{\text{lin}} = u \mathbf{r}_1 \mathbf{r}^\top = u \begin{bmatrix} \mathbf{r}_1^+ \\ \mathbf{r}_1^- \end{bmatrix} \mathbf{r}^\top, \quad (14a)$$

$$\mathbf{W}_l^{\text{lin}} = u \mathbf{r}_l \mathbf{r}_{l-1}^\top = u \begin{bmatrix} \mathbf{r}_l^+ \mathbf{r}_{l-1}^{+\top} & \mathbf{r}_l^+ \mathbf{r}_{l-1}^{-\top} \\ \mathbf{r}_l^- \mathbf{r}_{l-1}^{+\top} & \mathbf{r}_l^- \mathbf{r}_{l-1}^{-\top} \end{bmatrix}, \quad l = 2, \dots, L-1, \quad (14b)$$

$$\mathbf{W}_L^{\text{lin}} = u \mathbf{r}_{L-1}^\top = u \begin{bmatrix} \mathbf{r}_{L-1}^{+\top} & \mathbf{r}_{L-1}^{-\top} \end{bmatrix}, \quad (14c)$$

where u represents the norm of each layer, and $\mathbf{r}, \mathbf{r}_1, \mathbf{r}_2, \dots, \mathbf{r}_L$ are unit norm column vectors. The vectors $\mathbf{r}_l^+, \mathbf{r}_l^-$ denote the positive and negative elements in \mathbf{r}_l . The equal norm u of all layers is a consequence of small initialization (Du et al., 2018). Note that the weights can be written in blocks, as Equation (14), only after permuting the positive and negative elements. We use this permuted notation for the sake of exposition; no additional assumptions are required.

In a deep ReLU network, we empirically find that when the weights are trained from small initialization and the target function is linear, the weights form a particular rank-one and rank-two structure. The weights of a deep bias-free network can be written approximately as

$$\mathbf{W}_1 = u \mathbf{r}_1 \mathbf{r}^\top = u \begin{bmatrix} \mathbf{r}_1^+ \\ \mathbf{r}_1^- \end{bmatrix} \mathbf{r}^\top, \quad (15a)$$

$$\mathbf{W}_l = u \begin{bmatrix} \sqrt{2} \mathbf{r}_l^+ \mathbf{r}_{l-1}^{+\top} & \mathbf{0} \\ \mathbf{0} & \sqrt{2} \mathbf{r}_l^- \mathbf{r}_{l-1}^{-\top} \end{bmatrix}, \quad l = 2, \dots, L-1, \quad (15b)$$

$$\mathbf{W}_L = u \mathbf{r}_{L-1}^\top = u \begin{bmatrix} \mathbf{r}_{L-1}^{+\top} & \mathbf{r}_{L-1}^{-\top} \end{bmatrix}. \quad (15c)$$

For the first and last layers, the weights in the deep ReLU network have the same rank-one structure as their linear counterpart. For the intermediate layers, weights in the deep ReLU network are rank-two. Specifically, positive weights in the ReLU network correspond to positive weights in the linear network and zero weights in the ReLU network correspond to negative weights in the linear network. We visualize the low-rank weights of a three-layer linear network and a three-layer bias-free ReLU network in Figure 4.

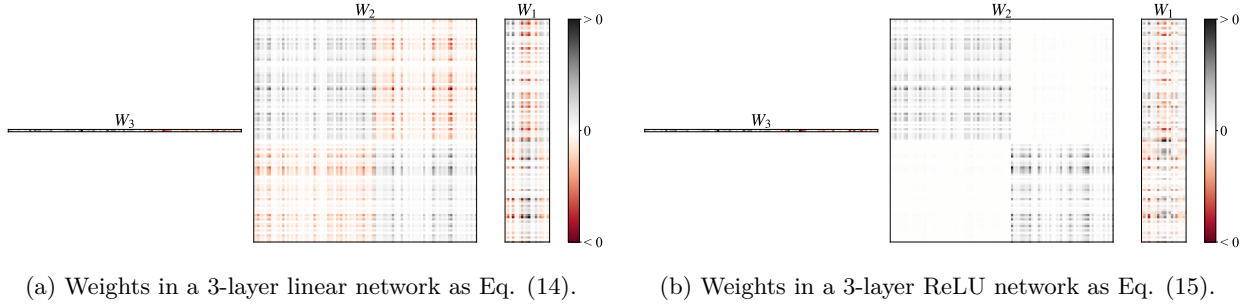


Figure 4: Low-rank weights in deep linear and ReLU bias-free networks. A three-layer linear network and a three-layer ReLU network are trained on the same dataset starting from the same small random weights. The dataset has a linear target function and an even empirical input data distribution. We plot the weights when the loss has approached zero. \mathbf{W}_1 , \mathbf{W}_3 , and positive elements in \mathbf{W}_2 have approximately the same structure in the linear and ReLU networks. Elements of \mathbf{W}_2 that are negative in the linear network are approximately zero in the ReLU network. The neurons are permuted for better visualization. Experimental details are provided in Appendix H.

If we assume $\|\mathbf{r}_l^+\| = \|\mathbf{r}_l^-\|$, ($l = 1, 2, \dots, L-1$), which is true when the network width approaches infinity, the deep bias-free ReLU network with weights defined in Equation (15) implements a linear function

$$f(\mathbf{x}; \mathbf{W}) = \mathbf{W}_L \mathbf{W}_{L-1} \cdots \mathbf{W}_2 \sigma(\mathbf{W}_1 \mathbf{x}) = \frac{1}{2} \mathbf{W}_L \cdots \mathbf{W}_2 \mathbf{W}_1 \mathbf{x}. \quad (16)$$

In the first equality, we drop the activation functions except the one between the first and second layers. This is because the second layer weights, \mathbf{W}_2 , is non-negative as shown in Figure 4b, and so is the output of a ReLU activation function, $\sigma(\mathbf{W}_1 \mathbf{x})$. Hence, their product, $\mathbf{W}_2 \sigma(\mathbf{W}_1 \mathbf{x})$, is also non-negative. We thus have $\sigma(\mathbf{W}_2 \sigma(\mathbf{W}_1 \mathbf{x})) = \mathbf{W}_2 \sigma(\mathbf{W}_1 \mathbf{x})$. The same applies to all subsequent layers. The second equality is obtained by substituting the weights defined in Equation (15) into the expression.

When the empirical input distribution is even and the target function is linear, the learning dynamics of the deep bias-free ReLU network with weights as Equation (15) reduces to (see Appendix E)

$$\tau \dot{\mathbf{W}}_l = \left(\prod_{l'=l+1}^L \mathbf{W}_{l'} \right)^\top \left(\frac{1}{2} \beta^\top - \frac{1}{4} \prod_{l'=1}^L \mathbf{W}_{l'} \Sigma \right) \left(\prod_{l'=l-1}^L \mathbf{W}_{l'} \right)^\top. \quad (17)$$

Except for constant coefficients, these equations are the same as that of the deep linear network given in Equation (7). We show, in Appendix E, that weights which have formed a low-rank structure as defined in Equation (15) maintain the structure over training.

We conjecture that in deep linear networks and certain deep ReLU networks, the weights of an intermediate layer align with the inputs to that layer. For example, the second-layer weight in a deep linear network is $\mathbf{W}_2^{\text{lin}} = u \mathbf{r}_2 \mathbf{r}_1^\top$ as given in Equation (14). Every row of $\mathbf{W}_2^{\text{lin}}$ aligns with \mathbf{r}_1^\top , which is parallel to any input to the second layer, $\mathbf{W}_1^{\text{lin}} \mathbf{x} = u \mathbf{r}_1 \mathbf{r}^\top \mathbf{x}$. The second-layer weight in the deep ReLU network is given in Equation (15). Some rows of \mathbf{W}_2 align with $\begin{bmatrix} \mathbf{r}_1^{+\top} & \mathbf{0} \end{bmatrix}$, which is parallel to some inputs ($\mathbf{r}^\top \mathbf{x} > 0$) to the second layer, $\sigma(\mathbf{W}_1 \mathbf{x}) = u \begin{bmatrix} \mathbf{r}_1^{+\top} \\ \mathbf{0} \end{bmatrix} \mathbf{r}^\top \mathbf{x}$. Other rows of \mathbf{W}_2 align with $\begin{bmatrix} \mathbf{0} & \mathbf{r}_1^{-\top} \end{bmatrix}$, which is parallel with inputs ($\mathbf{r}^\top \mathbf{x} < 0$) to the second layer, $\sigma(\mathbf{W}_1 \mathbf{x}) = u \begin{bmatrix} \mathbf{0} \\ \mathbf{r}_1^{-\top} \end{bmatrix} \mathbf{r}^\top \mathbf{x}$. This alignment phenomenon has been proven for deep linear networks (Ji & Telgarsky, 2019; Marion & Chizat, 2024). Here we empirically find similar phenomena in certain deep ReLU networks. Analytical proof of alignment in deep ReLU networks is an intriguing future direction.

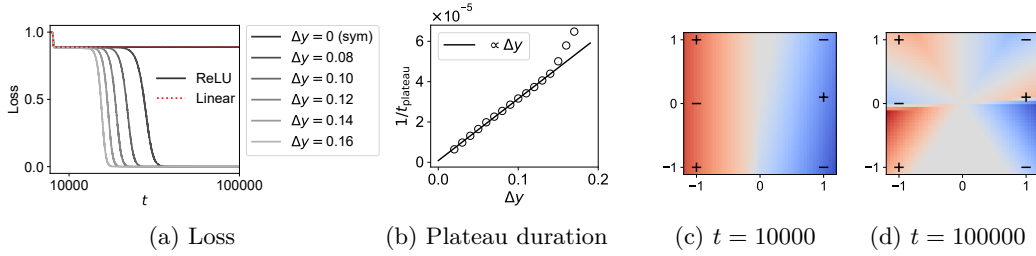


Figure 5: Two-layer bias-free linear/ReLU network trained on a dataset that slightly violates the symmetry Condition 3. The + and - signs represent data points with +1 and - labels respectively. The right middle data point is slight asymmetric with coordinates $(1, \Delta y)$. (a) Loss curves of the ReLU network with different Δy and the linear network with $\Delta y = 0.1$. (b) The duration of the plateau, during which the ReLU network implements a nearly linear solution, scales approximately with $1/\Delta y$. (c,d) The ReLU network output during and at the end of training. This ReLU network is trained on the dataset with $\Delta y = 0.1$. Experimental details are provided in Appendix H.

6 Discussion

Implication of Bias Removal. We studied the implications of removing bias in ReLU networks in terms of the expressivity and learning dynamics. Theorem 1 shows that two-layer bias-free (leaky) ReLU networks cannot express any odd functions except for linear functions. Theorem 7 shows that for datasets with an even input distribution and an odd target function, two-layer bias-free (leaky) ReLU networks have the same time evolution as a linear network (modulo scale factors) under initialization Assumption 5. We also presented examples in which the bias-free ReLU network evolves like multiple independent linear networks, in Section 4.2. In these cases, comparing a bias-free ReLU network with its linear counterpart provides an intuitive understanding of the behavior of ReLU networks. On the flip side, the simplicity of bias-free ReLU networks suggests that ReLU networks with bias may exhibit more complicated behaviors, which are not fully addressed by studies on bias-free networks, and remain open questions.

One common argument in studies of bias-free ReLU networks is that we can stack the input \mathbf{x} with an additional one, i.e., $[\mathbf{x}, 1]$. Then results derived for bias-free networks could extend to networks with bias and the removal of bias might thus have a minor implication. This argument is valid for some studies (Allen-Zhu et al., 2019; Zou et al., 2020), but not all. For example, Soudry et al. (2018) found that two-layer bias-free ReLU networks trained with logistic loss converge to the max-margin classifier on linearly separable datasets. As clarified by Soudry et al. (2018), this technical result holds when the inputs are stacked with an additional one. However, the max-margin solution for the dataset with stacked inputs is not the max-margin solution for the original dataset. Thus, the convergence to max-margin solution result does not directly extend to ReLU networks with bias.

Perturbed Symmetric Dataset. We have shown an exact equivalence between two-layer bias-free (leaky) ReLU networks and linear networks under symmetry Condition 3 on the dataset in Theorem 7. In practice, no datasets satisfies Condition 3 precisely. However, two-layer bias-free ReLU networks may still struggle to fit a dataset that approximately satisfies Condition 3. We present a simple example with six data points in Figure 5. The ReLU network loss curve closely follows the linear network loss curve in the early phase, when it first learns a nearly linear solution, as shown in Figures 5a and 5c. After a plateau, the ReLU network diverges from the linear network dynamics and converges to a nonlinear solution. The more symmetric the dataset, the longer the plateau a two-layer bias-free ReLU network undergoes before learning a nonlinear solution. As shown in Figure 5b, the inverse of the plateau duration scales approximately linearly with the deviation of the asymmetric data point, Δy . The scaling becomes less accurate for larger Δy because the corresponding dataset more severely violates symmetric Condition 3, where the ReLU network no longer behaves like a linear network. When the dataset is exactly symmetric, the ReLU network never learns a nonlinear solution. Furthermore, as shown in Figure 5d, the decision boundaries at convergence are close to the data points, and thus probably not robust.

References

- Madhu S. Advani, Andrew M. Saxe, and Haim Sompolsky. High-dimensional dynamics of generalization error in neural networks. *Neural Networks*, 132:428–446, 2020. ISSN 0893-6080. doi: <https://doi.org/10.1016/j.neunet.2020.08.022>. URL <https://www.sciencedirect.com/science/article/pii/S0893608020303117>.
- Zeyuan Allen-Zhu, Yuanzhi Li, and Zhao Song. A convergence theory for deep learning via overparameterization. In Kamalika Chaudhuri and Ruslan Salakhutdinov (eds.), *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pp. 242–252. PMLR, 09–15 Jun 2019. URL <https://proceedings.mlr.press/v97/allen-zhu19a.html>.
- Sanjeev Arora, Nadav Cohen, and Elad Hazan. On the optimization of deep networks: Implicit acceleration by overparameterization. In Jennifer Dy and Andreas Krause (eds.), *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pp. 244–253. PMLR, 10–15 Jul 2018. URL <https://proceedings.mlr.press/v80/arora18a.html>.
- Sanjeev Arora, Simon Du, Wei Hu, Zhiyuan Li, and Ruosong Wang. Fine-grained analysis of optimization and generalization for overparameterized two-layer neural networks. In Kamalika Chaudhuri and Ruslan Salakhutdinov (eds.), *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pp. 322–332. PMLR, 09–15 Jun 2019. URL <https://proceedings.mlr.press/v97/arora19a.html>.
- Alexander Atanasov, Blake Bordelon, and Cengiz Pehlevan. Neural networks as kernel learners: The silent alignment effect. In *International Conference on Learning Representations*, 2022. URL <https://openreview.net/forum?id=1NvflqAdoom>.
- Pierre Baldi and Kurt Hornik. Neural networks and principal component analysis: Learning from examples without local minima. *Neural Networks*, 2(1):53–58, 1989. ISSN 0893-6080. doi: [https://doi.org/10.1016/0893-6080\(89\)90014-2](https://doi.org/10.1016/0893-6080(89)90014-2). URL <https://www.sciencedirect.com/science/article/pii/0893608089900142>.
- Peter L. Bartlett, Philip M. Long, Gábor Lugosi, and Alexander Tsigler. Benign overfitting in linear regression. *Proceedings of the National Academy of Sciences*, 117(48):30063–30070, 2020. doi: 10.1073/pnas.1907378117. URL <https://www.pnas.org/doi/abs/10.1073/pnas.1907378117>.
- Ronen Basri, David Jacobs, Yoni Kasten, and Shira Kritchman. The convergence rate of neural networks for learned functions of different frequencies. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019. URL https://proceedings.neurips.cc/paper_files/paper/2019/file/5ac8bb8a7d745102a978c5f8ccdb61b8-Paper.pdf.
- Etienne Boursier and Nicolas Flammarion. Early alignment in two-layer networks training is a two-edged sword. *arXiv preprint arXiv:2401.10791*, 2024.
- Etienne Boursier, Loucas Pillaud-Vivien, and Nicolas Flammarion. Gradient flow dynamics of shallow relu networks for square loss and orthogonal inputs. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh (eds.), *Advances in Neural Information Processing Systems*, volume 35, pp. 20105–20118. Curran Associates, Inc., 2022. URL https://proceedings.neurips.cc/paper_files/paper/2022/file/7eeb9af3eb1f48e29c05e8dd3342b286-Paper-Conference.pdf.
- Lukas Braun, Clémentine Dominé, James Fitzgerald, and Andrew Saxe. Exact learning dynamics of deep linear networks with prior knowledge. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh (eds.), *Advances in Neural Information Processing Systems*, volume 35, pp. 6615–6629. Curran Associates, Inc., 2022. URL https://proceedings.neurips.cc/paper_files/paper/2022/file/2b3bb2c95195130977a51b3bb251c40a-Paper-Conference.pdf.

- Alon Brutzkus, Amir Globerson, Eran Malach, and Shai Shalev-Shwartz. SGD learns over-parameterized networks that provably generalize on linearly separable data. In *International Conference on Learning Representations*, 2018. URL <https://openreview.net/forum?id=rJ33wwxRb>.
- Dmitry Chistikov, Matthias Englert, and Ranko Lazic. Learning a neuron by a shallow relu network: Dynamics and implicit bias for correlated inputs. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine (eds.), *Advances in Neural Information Processing Systems*, volume 36, pp. 23748–23760. Curran Associates, Inc., 2023. URL https://proceedings.neurips.cc/paper_files/paper/2023/file/4af24e6ce753c181e703f3f0be3b5e20-Paper-Conference.pdf.
- Simon S Du, Wei Hu, and Jason D Lee. Algorithmic regularization in learning deep homogeneous models: Layers are automatically balanced. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 31. Curran Associates, Inc., 2018. URL https://proceedings.neurips.cc/paper_files/paper/2018/file/fe131d7f5a6b38b23cc967316c13dae2-Paper.pdf.
- Simon S. Du, Xiyu Zhai, Barnabas Poczos, and Aarti Singh. Gradient descent provably optimizes over-parameterized neural networks. In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=S1eK3i09YQ>.
- Spencer Frei, Niladri S. Chatterji, and Peter L. Bartlett. Random feature amplification: Feature learning and generalization in neural networks. *Journal of Machine Learning Research*, 24(303):1–49, 2023a. URL <http://jmlr.org/papers/v24/22-1132.html>.
- Spencer Frei, Gal Vardi, Peter Bartlett, and Nathan Srebro. Benign overfitting in linear classifiers and leaky relu networks from kkt conditions for margin maximization. In Gergely Neu and Lorenzo Rosasco (eds.), *Proceedings of Thirty Sixth Conference on Learning Theory*, volume 195 of *Proceedings of Machine Learning Research*, pp. 3173–3228. PMLR, 12–15 Jul 2023b. URL <https://proceedings.mlr.press/v195/frei23a.html>.
- Spencer Frei, Gal Vardi, Peter Bartlett, Nathan Srebro, and Wei Hu. Implicit bias in leaky reLU networks trained on high-dimensional data. In *The Eleventh International Conference on Learning Representations*, 2023c. URL <https://openreview.net/forum?id=JpbLyEI5EwW>.
- Spencer Frei, Gal Vardi, Peter Bartlett, and Nati Srebro. The double-edged sword of implicit bias: Generalization vs. robustness in relu networks. In A. Oh, T. Neumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine (eds.), *Advances in Neural Information Processing Systems*, volume 36, pp. 8885–8897. Curran Associates, Inc., 2023d. URL https://proceedings.neurips.cc/paper_files/paper/2023/file/1c26c389d60ec419fd24b5fee5b35796-Paper-Conference.pdf.
- Kenji Fukumizu. Effect of batch learning in multilayer neural networks. *Gen*, 1(04):1E–03, 1998.
- Gauthier Gidel, Francis Bach, and Simon Lacoste-Julien. Implicit regularization of discrete gradient dynamics in linear neural networks. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019. URL https://proceedings.neurips.cc/paper_files/paper/2019/file/f39ae9ff3a81f499230c4126e01f421b-Paper.pdf.
- Margalit Glasgow. SGD finds then tunes features in two-layer neural networks with near-optimal sample complexity: A case study in the XOR problem. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=Hg0JlxzB16>.
- T. H. Gronwall. Note on the derivatives with respect to a parameter of the solutions of a system of differential equations. *Annals of Mathematics*, 20(4):292–296, 1919. ISSN 0003486X. URL <http://www.jstor.org/stable/1967124>.
- David Holzmüller and Ingo Steinwart. Training two-layer relu networks with gradient descent is inconsistent. *Journal of Machine Learning Research*, 23(181):1–82, 2022. URL <http://jmlr.org/papers/v23/20-830.html>.

- Kurt Hornik, Maxwell Stinchcombe, and Halbert White. Multilayer feedforward networks are universal approximators. *Neural Networks*, 2(5):359–366, 1989. ISSN 0893-6080. doi: [https://doi.org/10.1016/0893-6080\(89\)90020-8](https://doi.org/10.1016/0893-6080(89)90020-8). URL <https://www.sciencedirect.com/science/article/pii/0893608089900208>.
- Arthur Jacot, François Ged, Berfin Şimşek, Clément Hongler, and Franck Gabriel. Saddle-to-saddle dynamics in deep linear networks: Small initialization training, symmetry, and sparsity. *arXiv preprint arXiv:2106.15933*, 2021.
- Ziwei Ji and Matus Telgarsky. Gradient descent aligns the layers of deep linear networks. In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=HJflg30qKX>.
- Zahra Kадkhodaie, Florentin Guth, Eero P Simoncelli, and Stéphane Mallat. Generalization in diffusion models arises from geometry-adaptive harmonic representations. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=ANvmVS2Yr0>.
- Yiwen Kou, Zixiang Chen, and Quanquan Gu. Implicit bias of gradient descent for two-layer relu and leaky relu networks on nearly-orthogonal data. In A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine (eds.), *Advances in Neural Information Processing Systems*, volume 36, pp. 30167–30221. Curran Associates, Inc., 2023. URL https://proceedings.neurips.cc/paper_files/paper/2023/file/602f5c1b803c53b2aaf0b3864bf3383a-Paper-Conference.pdf.
- Andrew K. Lampinen and Surya Ganguli. An analytic theory of generalization dynamics and transfer learning in deep linear networks. In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=ryfMLoCqtQ>.
- Thien Le and Stefanie Jegelka. Training invariances and the low-rank phenomenon: beyond linear networks. In *International Conference on Learning Representations*, 2022. URL <https://openreview.net/forum?id=XEW8CQgArno>.
- Tao Luo, Zhi-Qin John Xu, Zheng Ma, and Yaoyu Zhang. Phase diagram for two-layer relu neural networks at infinite-width limit. *Journal of Machine Learning Research*, 22(71):1–47, 2021. URL <http://jmlr.org/papers/v22/20-1123.html>.
- Kaifeng Lyu and Jian Li. Gradient descent maximizes the margin of homogeneous neural networks. In *International Conference on Learning Representations*, 2020. URL <https://openreview.net/forum?id=SJeLIgBKPS>.
- Kaifeng Lyu, Zhiyuan Li, Runzhe Wang, and Sanjeev Arora. Gradient descent on two-layer nets: Margin maximization and simplicity bias. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan (eds.), *Advances in Neural Information Processing Systems*, volume 34, pp. 12978–12991. Curran Associates, Inc., 2021. URL https://proceedings.neurips.cc/paper_files/paper/2021/file/6c351da15b5e8a743a21ee96a86e25df-Paper.pdf.
- Hartmut Maennel, Olivier Bousquet, and Sylvain Gelly. Gradient descent quantizes relu network features. *arXiv preprint arXiv:1803.08367*, 2018.
- Pierre Marion and Lénaïc Chizat. Deep linear networks for regression are implicitly regularized towards flat minima. In *Advances in Neural Information Processing Systems*, 2024. URL <https://openreview.net/forum?id=F738WY1Xm4>.
- Xuran Meng, Difan Zou, and Yuan Cao. Benign overfitting in two-layer relu convolutional neural networks for XOR data. In Ruslan Salakhutdinov, Zico Kolter, Katherine Heller, Adrian Weller, Nuria Oliver, Jonathan Scarlett, and Felix Berkenkamp (eds.), *Proceedings of the 41st International Conference on Machine Learning*, volume 235 of *Proceedings of Machine Learning Research*, pp. 35404–35469. PMLR, 21–27 Jul 2024. URL <https://proceedings.mlr.press/v235/meng24c.html>.

- Hancheng Min, Enrique Mallada, and Rene Vidal. Early neuron alignment in two-layer relu networks with small initialization. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=QibPzdVrRu>.
- Sreyas Mohan, Zahra Kadkhodaie, Eero P. Simoncelli, and Carlos Fernandez-Granda. Robust and interpretable blind image denoising via bias-free convolutional neural networks. In *International Conference on Learning Representations*, 2020. URL <https://openreview.net/forum?id=HJlSmC4FPS>.
- Mor Shpigel Nacson, Rotem Mulayoff, Greg Ongie, Tomer Michaeli, and Daniel Soudry. The implicit bias of minima stability in multivariate shallow reLU networks. In *The Eleventh International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=xtbog7cfsr>.
- Greg Ongie, Rebecca Willett, Daniel Soudry, and Nathan Srebro. A function space view of bounded norm infinite width relu nets: The multivariate case. In *International Conference on Learning Representations*, 2020. URL <https://openreview.net/forum?id=H1lNPxHKDH>.
- Scott Pesme and Nicolas Flammarion. Saddle-to-saddle dynamics in diagonal linear networks. In A. Oh, T. Neumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine (eds.), *Advances in Neural Information Processing Systems*, volume 36, pp. 7475–7505. Curran Associates, Inc., 2023. URL https://proceedings.neurips.cc/paper_files/paper/2023/file/17a9ab4190289f0e1504bbb98d1d111a-Paper-Conference.pdf.
- Leonardo Petrini, Francesco Cagnetta, Eric Vanden-Eijnden, and Matthieu Wyart. Learning sparse features can lead to overfitting in neural networks. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh (eds.), *Advances in Neural Information Processing Systems*, volume 35, pp. 9403–9416. Curran Associates, Inc., 2022. URL https://proceedings.neurips.cc/paper_files/paper/2022/file/3d3a9e085540c65dd3e5731361f9320e-Paper-Conference.pdf.
- Leonardo Petrini, Francesco Cagnetta, Eric Vanden-Eijnden, and Matthieu Wyart. Learning sparse features can lead to overfitting in neural networks. *Journal of Statistical Mechanics: Theory and Experiment*, 2023(11):114003, nov 2023. doi: 10.1088/1742-5468/ad01b9. URL <https://dx.doi.org/10.1088/1742-5468/ad01b9>.
- Mary Phuong and Christoph H Lampert. The inductive bias of relu networks on orthogonally separable data. In *International Conference on Learning Representations*, 2021. URL https://openreview.net/forum?id=krz7T0xU9Z_.
- Allan Pinkus. Approximation theory of the mlp model in neural networks. *Acta Numerica*, 8:143–195, 1999. doi: 10.1017/S0962492900002919.
- Roei Sarussi, Alon Brutzkus, and Amir Globerson. Towards understanding learning in neural networks with linear teachers. In Marina Meila and Tong Zhang (eds.), *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pp. 9313–9322. PMLR, 18–24 Jul 2021. URL <https://proceedings.mlr.press/v139/sarussi21a.html>.
- Andrew Saxe, Shagun Sodhani, and Sam Jay Lewallen. The neural race reduction: Dynamics of abstraction in gated networks. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvari, Gang Niu, and Sivan Sabato (eds.), *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pp. 19287–19309. PMLR, 17–23 Jul 2022. URL <https://proceedings.mlr.press/v162/saxe22a.html>.
- Andrew M Saxe, James L McClelland, and Surya Ganguli. Exact solutions to the nonlinear dynamics of learning in deep linear neural networks. In *International Conference on Learning Representations*, 2014. URL <https://arxiv.org/abs/1312.6120>.
- Andrew M. Saxe, James L. McClelland, and Surya Ganguli. A mathematical theory of semantic development in deep neural networks. *Proceedings of the National Academy of Sciences*, 116(23):11537–11546, 2019. doi: 10.1073/pnas.1820226116. URL <https://www.pnas.org/doi/abs/10.1073/pnas.1820226116>.

- Daniel Soudry, Elad Hoffer, Mor Shpigel Nacson, Suriya Gunasekar, and Nathan Srebro. The implicit bias of gradient descent on separable data. *Journal of Machine Learning Research*, 19(70):1–57, 2018. URL <http://jmlr.org/papers/v19/18-188.html>.
- Salma Tarmoun, Guilherme Franca, Benjamin D Haeffele, and Rene Vidal. Understanding the dynamics of gradient flow in overparameterized linear models. In Marina Meila and Tong Zhang (eds.), *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pp. 10153–10161. PMLR, 18–24 Jul 2021. URL <https://proceedings.mlr.press/v139/tarmoun21a.html>.
- Matus Telgarsky. Feature selection and low test error in shallow low-rotation relu networks. In *The Eleventh International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=swEskiem99>.
- Nadav Timor, Gal Vardi, and Ohad Shamir. Implicit regularization towards rank minimization in relu networks. In Shipra Agrawal and Francesco Orabona (eds.), *Proceedings of The 34th International Conference on Algorithmic Learning Theory*, volume 201 of *Proceedings of Machine Learning Research*, pp. 1429–1459. PMLR, 20 Feb–23 Feb 2023. URL <https://proceedings.mlr.press/v201/timor23a.html>.
- Gal Vardi and Ohad Shamir. Implicit regularization in relu networks with the square loss. In Mikhail Belkin and Samory Kpotufe (eds.), *Proceedings of Thirty Fourth Conference on Learning Theory*, volume 134 of *Proceedings of Machine Learning Research*, pp. 4224–4258. PMLR, 15–19 Aug 2021. URL <https://proceedings.mlr.press/v134/vardi21b.html>.
- Gal Vardi, Gilad Yehudai, and Ohad Shamir. Gradient methods provably converge to non-robust networks. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh (eds.), *Advances in Neural Information Processing Systems*, volume 35, pp. 20921–20932. Curran Associates, Inc., 2022. URL https://proceedings.neurips.cc/paper_files/paper/2022/file/83e6913572ba09b0ab53c64c016c7d1a-Paper-Conference.pdf.
- 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. doi: 10.1109/TSP.2019.2904921.
- Mingze Wang and Chao Ma. Understanding multi-phase optimization dynamics and rich nonlinear behaviors of reLU networks. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL <https://openreview.net/forum?id=konBXvt2iS>.
- Yifei Wang and Mert Pilanci. The convex geometry of backpropagation: Neural network gradient flows converge to extreme points of the dual convex program. In *International Conference on Learning Representations*, 2022. URL <https://openreview.net/forum?id=5QhUE1qiVC6>.
- Blake Woodworth, Suriya Gunasekar, Jason D. Lee, Edward Moroshko, Pedro Savarese, Itay Golan, Daniel Soudry, and Nathan Srebro. Kernel and rich regimes in overparametrized models. In Jacob Abernethy and Shivani Agarwal (eds.), *Proceedings of Thirty Third Conference on Learning Theory*, volume 125 of *Proceedings of Machine Learning Research*, pp. 3635–3673. PMLR, 09–12 Jul 2020. URL <https://proceedings.mlr.press/v125/woodworth20a.html>.
- Zhiwei Xu, Yutong Wang, Spencer Frei, Gal Vardi, and Wei Hu. Benign overfitting and grokking in relu networks for xor cluster data. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=BxHgpC6FNv>.
- John Zarka, Florentin Guth, and Stéphane Mallat. Separation and concentration in deep networks. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=8HhkbjrWLdE>.
- Kai Zhang, Yawei Li, Wangmeng Zuo, Lei Zhang, Luc Van Gool, and Radu Timofte. Plug-and-play image restoration with deep denoiser prior. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(10):6360–6376, 2022. doi: 10.1109/TPAMI.2021.3088914.

- Xiao Zhang, Yaodong Yu, Lingxiao Wang, and Quanquan Gu. Learning one-hidden-layer relu networks via gradient descent. In Kamalika Chaudhuri and Masashi Sugiyama (eds.), *Proceedings of the Twenty-Second International Conference on Artificial Intelligence and Statistics*, volume 89 of *Proceedings of Machine Learning Research*, pp. 1524–1534. PMLR, 16–18 Apr 2019. URL <https://proceedings.mlr.press/v89/zhang19g.html>.
- Yedi Zhang, Peter E. Latham, and Andrew M Saxe. Understanding unimodal bias in multimodal deep linear networks. In Ruslan Salakhutdinov, Zico Kolter, Katherine Heller, Adrian Weller, Nuria Oliver, Jonathan Scarlett, and Felix Berkenkamp (eds.), *Proceedings of the 41st International Conference on Machine Learning*, volume 235 of *Proceedings of Machine Learning Research*, pp. 59100–59125. PMLR, 21–27 Jul 2024. URL <https://proceedings.mlr.press/v235/zhang24aa.html>.
- Liu Ziyin, Botao Li, and Xiangming Meng. Exact solutions of a deep linear network. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh (eds.), *Advances in Neural Information Processing Systems*, volume 35, pp. 24446–24458. Curran Associates, Inc., 2022. URL https://proceedings.neurips.cc/paper_files/paper/2022/file/9a940e858b17f01c402e164835140c4a-Paper-Conference.pdf.
- Difan Zou, Yuan Cao, Dongruo Zhou, and Quanquan Gu. Gradient descent optimizes over-parameterized deep relu networks. *Machine learning*, 109:467–492, 2020.

A Additional Related Work

Implicit / Simplicity Bias. Many works have studied the implicit bias or simplicity bias of two-layer bias-free ReLU networks under various assumptions on the dataset. Brutzkus et al. (2018); Wang et al. (2019); Lyu et al. (2021); Sarussi et al. (2021); Wang & Ma (2023) considered linearly separable binary classification tasks. Phuong & Lampert (2021); Wang & Pilanci (2022); Min et al. (2024) studied orthogonally separable classification (i.e., where for every pair of labeled examples $(\mathbf{x}_i, y_i), (\mathbf{x}_j, y_j)$ we have $\mathbf{x}_i^\top \mathbf{x}_j > 0$ if $y_i = y_j$ and $\mathbf{x}_i^\top \mathbf{x}_j \leq 0$ if otherwise). Boursier et al. (2022); Frei et al. (2023b;c); Kou et al. (2023) studied binary classification with exactly or nearly orthogonal input (i.e., where $\mathbf{x}_i^\top \mathbf{x}_j = 0$ if $i \neq j$). Orthogonal input is a sufficient condition for linear separability for binary classification tasks. Frei et al. (2023a); Meng et al. (2024); Xu et al. (2024) studied XOR-like datasets. Vardi et al. (2022); Frei et al. (2023d) studied datasets with adversarial noise. We add to this line of research by studying a case with extreme simplicity bias, i.e., behaving like linear networks.

Low-Rank Weights. Maennel et al. (2018) is the seminal work on the low-rank weights in two-layer ReLU networks trained from small initialization. They described the phenomenon as “quantizing”, where the first layer weight vectors align with a small number of directions in the early phase of training. Luo et al. (2021) identified when two-layer bias-free ReLU networks form low-rank weights in terms of the initialization and the network width. Timor et al. (2023) provided cases where gradient flow on two-layer and deep ReLU networks provably minimize or not minimize the ranks of weight matrices. Frei et al. (2023c); Kou et al. (2023) computed the numerical rank of the converged weights of two-layer bias-free ReLU networks for nearly orthogonal datasets, and found that weights in leaky ReLU networks have rank at most two and weights in ReLU networks have a numerical rank upper bounded by a constant. Chistikov et al. (2023) showed two-layer bias-free ReLU networks are implicitly biased to learn the network of minimal rank under the assumption that training points are correlated with the teacher neuron. Min et al. (2024); Boursier & Flammarion (2024) studied the early phase learning dynamics to understand how the low-rank weights form. Petrini et al. (2022; 2023) conducted experiments on practical datasets to show that two-layer bias-free ReLU networks learn sparse features, which can be detrimental and lead to overfitting. Le & Jegelka (2022) generalize the low-rank phenomenon in linear and ReLU networks to arbitrary non-homogeneous networks whose last few layers contain linear fully-connected and linear ResNet blocks.

B Useful Lemmas

Grönwall’s Inequality (Gronwall, 1919) is a common tool to obtain error bounds when considering approximate differential equations.

Lemma 10 (Grönwall’s Inequality). *Let I denote an interval of the real line of the form $[a, \infty)$ or $[a, b]$ or $[a, b]$ with $a < b$. Let α, β and u be real-valued functions defined on I . Assume that β and u are continuous and that the negative part of α is integrable on every closed and bounded subinterval of I . If β is non-negative and u satisfies the integral inequality*

$$u(t) \leq \alpha(t) + \int_a^t \beta(s)u(s)ds, \quad \forall t \in I,$$

then

$$u(t) \leq \alpha(t) + \int_a^t \alpha(s)\beta(s)e^{\int_s^t \beta(r)dr}ds, \quad t \in I. \quad (18)$$

The key implication of Condition 3 we exploit is that the input covariance matrix and the input-output correlation averaged over any half space is equal to those averaged over the entire space.

Lemma 11. *Let set \mathbb{S}^+ be an arbitrary half space divided by a hyperplane with normal vector \mathbf{r} , namely $\mathbb{S}^+ = \{\mathbf{x} \in \mathbb{R}^D | \mathbf{r}^\top \mathbf{x} > 0\}$. Under the first condition in Condition 3, we have $\forall \mathbf{r}$*

$$\langle \mathbf{x}\mathbf{x}^\top \rangle_{\mathbb{S}^+} \equiv \frac{\int_{\mathbb{S}^+} \mathbf{x}\mathbf{x}^\top p(\mathbf{x})d\mathbf{x}}{\int_{\mathbb{S}^+} p(\mathbf{x})d\mathbf{x}} = \Sigma. \quad (19)$$

Under Condition 3, we have $\forall \mathbf{r}$

$$\langle \mathbf{x}y(\mathbf{x}) \rangle_{\mathbb{S}^+} \equiv \frac{\int_{\mathbb{S}^+} \mathbf{x}y(\mathbf{x})p(\mathbf{x})d\mathbf{x}}{\int_{\mathbb{S}^+} p(\mathbf{x})d\mathbf{x}} = \beta \quad (20)$$

Recall that Σ and β are averages over the entire space as defined in Equation (4).

Proof. Define $\mathbb{S}^- = \{\mathbf{x} \in \mathbb{R}^D | \mathbf{r}^\top \mathbf{x} < 0\}$. Because Condition 3 states that $p(\mathbf{x})$ is even, we have

$$\int_{\mathbb{S}^+} p(\mathbf{x})d\mathbf{x} = \int_{\mathbb{S}^-} p(\mathbf{x})d\mathbf{x} = \frac{1}{2}.$$

Because $\mathbf{x}\mathbf{x}^\top p(\mathbf{x})$ is an even function about \mathbf{x} , the integral of $\mathbf{x}\mathbf{x}^\top p(\mathbf{x})$ over \mathbb{S}^+ or \mathbb{S}^- is the same

$$\int_{\mathbb{S}^+} \mathbf{x}\mathbf{x}^\top p(\mathbf{x})d\mathbf{x} = \int_{\mathbb{S}^-} \mathbf{x}\mathbf{x}^\top p(\mathbf{x})d\mathbf{x}.$$

Thus the average of $\mathbf{x}\mathbf{x}^\top$ over \mathbb{S}^+ is equal to the average in the entire space,

$$\Sigma = \int_{\mathbb{S}^+} \mathbf{x}\mathbf{x}^\top p(\mathbf{x})d\mathbf{x} + \int_{\mathbb{S}^-} \mathbf{x}\mathbf{x}^\top p(\mathbf{x})d\mathbf{x} = 2 \int_{\mathbb{S}^+} \mathbf{x}\mathbf{x}^\top p(\mathbf{x})d\mathbf{x} = \langle \mathbf{x}\mathbf{x}^\top \rangle_{\mathbb{S}^+}.$$

The same holds for $\mathbf{x}y(\mathbf{x})$. □

Lemma 12. Under Condition 3, the first terms in the differential Equation (2) reduce to

$$\langle \sigma'(\mathbf{W}_1 \mathbf{x}) \odot \mathbf{W}_2^\top y \mathbf{x}^\top \rangle = \frac{\alpha + 1}{2} \mathbf{W}_2^\top \beta^\top, \quad (21a)$$

$$\langle y \sigma(\mathbf{W}_1 \mathbf{x})^\top \rangle = \frac{\alpha + 1}{2} \beta^\top \mathbf{W}_1^\top. \quad (21b)$$

Proof. Let us consider the h -th row of the matrix $\langle \sigma'(\mathbf{W}_1 \mathbf{x}) \odot \mathbf{W}_2^\top y \mathbf{x}^\top \rangle$, which is

$$\langle \sigma'(\mathbf{w}_{1h} \mathbf{x}) w_{2h} y \mathbf{x}^\top \rangle = \frac{1}{2} \langle \alpha w_{2h} y \mathbf{x}^\top \rangle_{\mathbf{w}_{1h} \mathbf{x} < 0} + \frac{1}{2} \langle w_{2h} y \mathbf{x}^\top \rangle_{\mathbf{w}_{1h} \mathbf{x} > 0} = \frac{\alpha + 1}{2} w_{2h} \beta^\top,$$

where the first equality is the law of total expectation and the second equality uses Equation (20). Because the same holds for all rows, Equation (21a) is true.

Let us consider the h -th element of the row vector $\langle y \sigma(\mathbf{W}_1 \mathbf{x})^\top \rangle$, which is

$$\langle y \sigma(\mathbf{w}_{1h} \mathbf{x}) \rangle = \frac{1}{2} \langle \alpha y \mathbf{w}_{1h} \mathbf{x} \rangle_{\mathbf{w}_{1h} \mathbf{x} < 0} + \frac{1}{2} \langle y \mathbf{w}_{1h} \mathbf{x} \rangle_{\mathbf{w}_{1h} \mathbf{x} > 0} = \frac{\alpha + 1}{2} \mathbf{w}_{1h} \beta,$$

where the first equality is the law of total expectation and the second equality uses Equation (20). Because the same holds for all elements, Equation (21b) is true. □

Lemma 13. The second terms in the differential Equation (2) can be bounded by the norm of the weights and the trace of the input covariance matrix.

1. $\|\langle \sigma(\mathbf{W}_1 \mathbf{x}) \sigma(\mathbf{W}_1 \mathbf{x})^\top \rangle\| \leq \|\mathbf{W}_1\|^2 \text{Tr } \Sigma.$
2. $\|\langle \sigma'(\mathbf{W}_1 \mathbf{x}) \odot \mathbf{W}_2^\top \mathbf{W}_2 \sigma(\mathbf{W}_1 \mathbf{x}) \mathbf{x}^\top \rangle\| \leq \|\mathbf{W}_2\|^2 \|\mathbf{W}_1\| \text{Tr } \Sigma.$

Proof. For the first inequality,

$$\begin{aligned} \|\langle \sigma(\mathbf{W}_1 \mathbf{x}) \sigma(\mathbf{W}_1 \mathbf{x})^\top \rangle\| &\leq \langle \|\sigma(\mathbf{W}_1 \mathbf{x})\|^2 \rangle \\ &\leq \langle \|\mathbf{W}_1 \mathbf{x}\|^2 \rangle \\ &\leq \langle \|\mathbf{W}_1\|^2 \|\mathbf{x}\|^2 \rangle \\ &= \|\mathbf{W}_1\|^2 \text{Tr } \Sigma. \end{aligned}$$

For the second inequality,

$$\begin{aligned}
\|\langle \sigma'(\mathbf{W}_1 \mathbf{x}) \odot \mathbf{W}_2^\top \mathbf{W}_2 \sigma(\mathbf{W}_1 \mathbf{x}) \mathbf{x}^\top \rangle\| &\leq \|\langle \mathbf{W}_2^\top \mathbf{W}_2 \sigma(\mathbf{W}_1 \mathbf{x}) \mathbf{x}^\top \rangle\| \\
&\leq \|\mathbf{W}_2\|^2 \langle \|\mathbf{W}_1 \mathbf{x}\| \|\mathbf{x}\| \rangle \\
&\leq \|\mathbf{W}_2\|^2 \langle \|\mathbf{W}_1\| \|\mathbf{x}\|^2 \rangle \\
&= \|\mathbf{W}_2\|^2 \|\mathbf{W}_1\| \text{Tr } \Sigma.
\end{aligned}$$

□

C Two-Layer Networks on Symmetric Datasets

C.1 Learning Dynamics: Early Phase

In the early phase of learning, the network output is small compared to the target output because the initialization is small. Lemma 14 specifies how small the norm of the weights is.

Lemma 14. *Denote the larger L2 norm of the weights in a two-layer network as $u(t) = \max\{\|\mathbf{W}_1(t)\|, \|\mathbf{W}_2(t)\|\}$. The initial weights are small, that is $w_{\text{init}} \equiv u(0) \ll 1$. For two-layer linear, ReLU, or leaky ReLU networks trained with square loss from small initialization, $u(t)$ is bounded by*

$$u(t) \leq u(0)e^{(s+\text{Tr } \Sigma)t/\tau}, \quad (22)$$

for time $t < \frac{\tau}{s+\text{Tr } \Sigma} \ln \frac{1}{w_{\text{init}}}$.

Proof. For two-layer linear, ReLU, or leaky ReLU networks, the learning dynamics are given in general in Equation (2). Using the equality in Lemma 12 and the inequality in Lemma 13, we can bound the dynamics of u^2 as

$$\begin{aligned}
\tau \frac{d}{dt} u^2 &= \tau \frac{d}{dt} \|\mathbf{W}_2\|^2 = (\alpha + 1) \beta^\top \mathbf{W}_1^\top \mathbf{W}_2^\top - 2 \mathbf{W}_2 \langle \sigma(\mathbf{W}_1 \mathbf{x}) \sigma(\mathbf{W}_1 \mathbf{x})^\top \rangle \mathbf{W}_2^\top \\
&\leq |(\alpha + 1) \beta^\top \mathbf{W}_1^\top \mathbf{W}_2^\top| + |2 \mathbf{W}_2 \langle \sigma(\mathbf{W}_1 \mathbf{x}) \sigma(\mathbf{W}_1 \mathbf{x})^\top \rangle \mathbf{W}_2^\top| \\
&\leq 2 \|\beta\| \|\mathbf{W}_1\| \|\mathbf{W}_2\| + 2 \|\mathbf{W}_2\|^2 \|\mathbf{W}_1\|^2 \text{Tr } \Sigma \\
&\leq 2su^2 + 2u^4 \text{Tr } \Sigma.
\end{aligned}$$

where $s = \|\beta\|$. For $u < 1$, we have

$$\tau \frac{d}{dt} u^2 \leq 2su^2 + 2u^4 \text{Tr } \Sigma < 2(s + \text{Tr } \Sigma) u^2.$$

Via Lemma 10 Grönwall's Inequality, we obtain

$$u^2 \leq w_{\text{init}}^2 e^{2(s+\text{Tr } \Sigma)t/\tau} \Rightarrow u(t) \leq w_{\text{init}} e^{(s+\text{Tr } \Sigma)t/\tau}.$$

This holds for

$$t < \frac{\tau}{s + \text{Tr } \Sigma} \ln \frac{1}{w_{\text{init}}}.$$

□

Since the weights are small in the early phase, we can approximate the early phase dynamics with only the first terms in Equation (2), that is

$$\tau \dot{\mathbf{W}}_1 \approx \langle \sigma'(\mathbf{W}_1 \mathbf{x}) \odot \mathbf{W}_2^\top \mathbf{y} \mathbf{x}^\top \rangle = \frac{\alpha + 1}{2} \mathbf{W}_2^\top \beta^\top, \quad (23)$$

$$\tau \dot{\mathbf{W}}_2 \approx \langle \mathbf{y} \sigma(\mathbf{W}_1 \mathbf{x})^\top \rangle = \frac{\alpha + 1}{2} \beta^\top \mathbf{W}_1^\top, \quad (24)$$

where the equalities hold under Condition 3 as proven by Lemma 12. We solve the approximate early phase dynamics and prove that the errors introduced by the approximation are bounded in Theorem 15.

Theorem 15. For time $t < \frac{\tau}{s+\text{Tr } \Sigma} \ln \frac{1}{w_{\text{init}}}$, the solution to Equation (2) starting from small initialization is exponential growth along one direction with small errors

$$\mathbf{W}_1(t) = e^{\frac{\alpha+1}{2\tau} st} \mathbf{r}_1 \bar{\beta}^\top + O(w_{\text{init}}), \quad \mathbf{W}_2(t) = e^{\frac{\alpha+1}{2\tau} st} \mathbf{r}_1^\top + O(w_{\text{init}}). \quad (25)$$

where $s = \|\beta\|$, $\bar{\beta} = \beta/s$, and \mathbf{r}_1 is determined by random initialization $\mathbf{r}_1 = (\mathbf{W}_1(0)\bar{\beta} + \mathbf{W}_2^\top(0))/2$.

Proof. We first consider the approximate learning dynamics:

$$\tau \dot{\widetilde{\mathbf{W}}}_1 = \frac{\alpha+1}{2} \widetilde{\mathbf{W}}_2^\top \beta^\top, \quad \tau \dot{\widetilde{\mathbf{W}}}_2 = \frac{\alpha+1}{2} \beta^\top \widetilde{\mathbf{W}}_1^\top. \quad (26)$$

This is a linear dynamical system with an analytical solution available. We re-write it as:

$$\tau \frac{d}{dt} \widetilde{\mathbf{W}} = \frac{\alpha+1}{2} \mathbf{M} \widetilde{\mathbf{W}}, \quad \text{where } \mathbf{M} = \begin{bmatrix} \mathbf{0} & \beta \\ \beta^\top & 0 \end{bmatrix}, \quad \widetilde{\mathbf{W}} = \begin{bmatrix} \widetilde{\mathbf{W}}_1^\top \\ \widetilde{\mathbf{W}}_2 \end{bmatrix}. \quad (27)$$

Since matrix \mathbf{M} only has two nonzero eigenvalues $\pm s$, the solution to Equation (27) is

$$\begin{aligned} \widetilde{\mathbf{W}}(t) &= \frac{1}{2} e^{\frac{\alpha+1}{2\tau} st} \begin{bmatrix} \bar{\beta} \\ 1 \end{bmatrix} (\bar{\beta}^\top \mathbf{W}_1^\top(0) + \mathbf{W}_2(0)) \\ &\quad + \frac{1}{2} e^{-\frac{\alpha+1}{2\tau} st} \begin{bmatrix} \bar{\beta} \\ -1 \end{bmatrix} (\bar{\beta}^\top \mathbf{W}_1^\top(0) - \mathbf{W}_2(0)) + \begin{bmatrix} (\mathbf{I} - \bar{\beta} \bar{\beta}^\top) \mathbf{W}_1^\top(0) \\ 0 \end{bmatrix}. \end{aligned} \quad (28)$$

Note that only the first term in Equation (28) is growing.

We then consider the exact learning dynamics given by Equation (2) and prove its solution is close to $\widetilde{\mathbf{W}}(t)$. The dynamics of the difference between the exact and approximate dynamics are

$$\tau \frac{d}{dt} (\widetilde{\mathbf{W}}_1 - \mathbf{W}_1) = \frac{\alpha+1}{2} (\widetilde{\mathbf{W}}_2 - \mathbf{W}_2)^\top \beta^\top + \langle \sigma'(\mathbf{W}_1 \mathbf{x}) \odot \mathbf{W}_2^\top \mathbf{W}_2 \sigma(\mathbf{W}_1 \mathbf{x}) \mathbf{x}^\top \rangle \quad (29a)$$

$$\tau \frac{d}{dt} (\widetilde{\mathbf{W}}_2 - \mathbf{W}_2) = \frac{\alpha+1}{2} \beta^\top (\widetilde{\mathbf{W}}_1 - \mathbf{W}_1)^\top + \mathbf{W}_2 \langle \sigma(\mathbf{W}_1 \mathbf{x}) \sigma(\mathbf{W}_1 \mathbf{x})^\top \rangle. \quad (29b)$$

We re-write Equation (29) as

$$\tau \frac{d}{dt} \delta \mathbf{W} = \frac{\alpha+1}{2} \mathbf{M} \delta \mathbf{W} + \epsilon, \quad (30)$$

The norm of the two components of ϵ can be bounded via Lemma 13

$$\begin{aligned} \|\langle \sigma'(\mathbf{W}_1 \mathbf{x}) \odot \mathbf{W}_2^\top \mathbf{W}_2 \sigma(\mathbf{W}_1 \mathbf{x}) \mathbf{x}^\top \rangle\| &\leq u^3 \text{Tr } \Sigma, \\ \|\mathbf{W}_2 \langle \sigma(\mathbf{W}_1 \mathbf{x}) \sigma(\mathbf{W}_1 \mathbf{x})^\top \rangle\| &\leq u^3 \text{Tr } \Sigma. \end{aligned}$$

We can then substitute in Equation (22) and obtain

$$\|\epsilon\| \leq \sqrt{2} u^3 \text{Tr } \Sigma < \sqrt{2} u_0^3 e^{3(s+\text{Tr } \Sigma)t/\tau} \text{Tr } \Sigma.$$

We now bound the norm of $\mathbf{W} - \widetilde{\mathbf{W}}$

$$\begin{aligned} \|\mathbf{W} - \widetilde{\mathbf{W}}\| &= \left\| \int_0^t \frac{\alpha+1}{2} \mathbf{M} (\mathbf{W} - \widetilde{\mathbf{W}}) + \epsilon dt \right\| \\ &\leq \int_0^t \|\mathbf{M}\| \|\mathbf{W} - \widetilde{\mathbf{W}}\| + \|\epsilon\| dt \\ &\leq \int_0^t \left(\sqrt{2} s \|\mathbf{W} - \widetilde{\mathbf{W}}\| + \sqrt{2} u_0^3 e^{3(s+\text{Tr } \Sigma)t/\tau} \text{Tr } \Sigma \right) dt \\ &\leq \frac{\sqrt{2} u_0^3 \text{Tr } \Sigma}{3(s+\text{Tr } \Sigma)} \left(e^{3(s+\text{Tr } \Sigma)t/\tau} - 1 \right) + \sqrt{2} s \int_0^t \|\mathbf{W} - \widetilde{\mathbf{W}}\| dt. \end{aligned}$$

Via Lemma 10 Grönwall's Inequality, we obtain

$$\begin{aligned}\|\mathbf{W} - \widetilde{\mathbf{W}}\| &\leq \frac{\sqrt{2} \operatorname{Tr} \Sigma u_0^3}{3(s + \operatorname{Tr} \Sigma)} \left[e^{3(s + \operatorname{Tr} \Sigma)t/\tau} - 1 + \int_0^t \left(e^{3(s + \operatorname{Tr} \Sigma)t'/\tau} - 1 \right) \sqrt{2} s e^{\sqrt{2} s t'} dt' \right] \\ &= \frac{\sqrt{2} \operatorname{Tr} \Sigma u_0^3}{3(s + \operatorname{Tr} \Sigma)} \left[e^{3(s + \operatorname{Tr} \Sigma)t/\tau} + \frac{\sqrt{2} s \left(e^{[3(s + \operatorname{Tr} \Sigma) + \sqrt{2} s]t/\tau} - 1 \right)}{3(s + \operatorname{Tr} \Sigma) + \sqrt{2} s} - e^{\sqrt{2} s t} \right].\end{aligned}$$

When $t < \frac{\tau}{s + \operatorname{Tr} \Sigma} \ln \frac{1}{u_0}$, we have $\|\mathbf{W} - \widetilde{\mathbf{W}}\| < C_1 u_0^2$ for some constant C_1 .

We are now ready to bound the difference between the exact solution and an exponential function along one direction

$$\begin{aligned}\mathbf{W} - e^{\frac{\alpha+1}{2\tau} s t} \begin{bmatrix} \bar{\beta} \\ 1 \end{bmatrix} \mathbf{r}_1^\top \\ &= (\mathbf{W} - \widetilde{\mathbf{W}}) + \left(\widetilde{\mathbf{W}} - e^{-\frac{\alpha+1}{2\tau} s t} \begin{bmatrix} \bar{\beta} \\ 1 \end{bmatrix} \mathbf{r}_1^\top \right) \\ &= (\mathbf{W} - \widetilde{\mathbf{W}}) + \frac{1}{2} e^{-\frac{\alpha+1}{2\tau} s t} \begin{bmatrix} \bar{\beta} \\ -1 \end{bmatrix} (\bar{\beta}^\top \mathbf{W}_1^\top(0) - \mathbf{W}_2(0)) + \begin{bmatrix} (I - \bar{\beta} \bar{\beta}^\top) \mathbf{W}_1^\top(0) \\ 0 \end{bmatrix}.\end{aligned}$$

The first term arises from our approximation of dropping the cubic terms in the dynamics. Its norm is bounded by $C_1 w_{\text{init}}^2$. The second term arises from initialization, which is $O(w_{\text{init}})$. Via triangle inequality, the norm of the total error is of order $O(w_{\text{init}})$.

$$\left\| \mathbf{W} - e^{\frac{\alpha+1}{2\tau} s t} \begin{bmatrix} \bar{\beta} \\ 1 \end{bmatrix} \mathbf{r}_1^\top \right\| < C_1 w_{\text{init}}^2 + C_2 w_{\text{init}} < C w_{\text{init}}.$$

□

Theorem 15 implies two messages. Firstly, the (leaky) ReLU network has the same time-course solution as its linear counterpart except a scale factor determined by α , which is consistent with Theorem 7. Secondly, the (leaky) ReLU and linear networks form rank-one weights with small errors in the early phase. We exploit the rank-one weights to reduce the learning dynamics to Equation (10).

C.2 Learning Dynamics: Late Phase

Proof of Theorem 7 (square loss). Theorem 7 relies on Condition 3 and Assumption 5 and arrives at three statements: implementing the same function as in Equation (11), having the same weights as in Equation (12), and retaining rank-one weights as Assumption 5. We prove them one by one.

Part 1: We first prove that the (leaky) ReLU network and the linear network implement the same linear function except scaling when their weights satisfy Assumption 5. Denote $\mathbf{W}_2 = [\mathbf{W}_2^+, \mathbf{W}_2^-]$ where \mathbf{W}_2^+ are the positive elements in \mathbf{W}_2 and \mathbf{W}_2^- are the negative elements in \mathbf{W}_2 . For a (leaky) ReLU network with rank-one weights satisfying Assumption 5, we have

$$f(\mathbf{x}; \mathbf{W}) = \mathbf{W}_2 \sigma(\mathbf{W}_1 \mathbf{x}) = \mathbf{W}_2 \sigma(\mathbf{W}_2^\top \mathbf{r}^\top \mathbf{x}).$$

Notate the positive and negative half-space as

$$\mathbb{S}^+ = \{\mathbf{x} \in \mathbb{R}^D \mid \mathbf{r}^\top \mathbf{x} \geq 0\}, \quad \mathbb{S}^- = \{\mathbf{x} \in \mathbb{R}^D \mid \mathbf{r}^\top \mathbf{x} < 0\}. \quad (31)$$

For $\mathbf{x} \in \mathbb{S}^+$, we have

$$f(\mathbf{x}; \mathbf{W}) = \mathbf{r}^\top \mathbf{x} \mathbf{W}_2 \sigma(\mathbf{W}_2^\top) = \mathbf{r}^\top \mathbf{x} \begin{bmatrix} \mathbf{W}_2^+ & \mathbf{W}_2^- \end{bmatrix} \begin{bmatrix} \mathbf{W}_2^{+\top} \\ \alpha \mathbf{W}_2^{-\top} \end{bmatrix} = \mathbf{r}^\top \mathbf{x} (\|\mathbf{W}_2^+\|^2 + \alpha \|\mathbf{W}_2^-\|^2).$$

For $\mathbf{x} \in \mathbb{S}^-$, we have

$$f(\mathbf{x}; \mathbf{W}) = -\mathbf{r}^\top \mathbf{x} \mathbf{W}_2 \sigma(-\mathbf{W}_2^\top) = \mathbf{r}^\top \mathbf{x} \begin{bmatrix} \mathbf{W}_2^+ & \mathbf{W}_2^- \end{bmatrix} \begin{bmatrix} \alpha \mathbf{W}_2^{+\top} \\ \mathbf{W}_2^{-\top} \end{bmatrix} = \mathbf{r}^\top \mathbf{x} (\alpha \|\mathbf{W}_2^+\|^2 + \|\mathbf{W}_2^-\|^2).$$

Under Assumption 5, we have $\|\mathbf{W}_2^+\| = \|\mathbf{W}_2^-\|$. Hence, the (leaky) ReLU network implements

$$\forall \mathbf{r}, \quad f(\mathbf{x}; \mathbf{W}) = \mathbf{W}_2 \sigma(\mathbf{W}_1 \mathbf{x}) = \frac{\alpha+1}{2} \mathbf{r}^\top \mathbf{x} \|\mathbf{W}_2\|^2. \quad (32)$$

Under Assumption 5, the linear network implements

$$f^{\text{lin}}(\mathbf{x}; \mathbf{W}) = \mathbf{W}_2 \mathbf{W}_1 \mathbf{x} = \mathbf{W}_2 \mathbf{W}_2^\top \mathbf{r}^\top \mathbf{x} = \mathbf{r}^\top \mathbf{x} \|\mathbf{W}_2\|^2. \quad (33)$$

Comparing Equations (32) and (33), we find that when the weights satisfy Assumption 5, the (leaky) ReLU network implements the same function as the linear network except a scale factor

$$\mathbf{W}_2 \sigma(\mathbf{W}_1 \mathbf{x}) = \frac{\alpha+1}{2} \mathbf{W}_2 \mathbf{W}_1 \mathbf{x}. \quad (34)$$

Part 2: We then look into the learning dynamics to prove that the weights in the (leaky) ReLU and the linear network are the same except scaling. Substituting Equation (34) into the dynamics, we get

$$\begin{aligned} \tau \dot{\mathbf{W}}_1 &= \frac{\alpha+1}{2} \mathbf{W}_2^\top \beta^\top - \langle \sigma'(\mathbf{W}_1 \mathbf{x}) \odot \mathbf{W}_2^\top \mathbf{W}_2 \sigma(\mathbf{W}_1 \mathbf{x}) \mathbf{x}^\top \rangle \\ &= \frac{\alpha+1}{2} \mathbf{W}_2^\top \beta^\top - \frac{\alpha+1}{2} \langle \sigma'(\mathbf{W}_1 \mathbf{x}) \odot \mathbf{W}_2^\top \mathbf{W}_2 \mathbf{W}_1 \mathbf{x} \mathbf{x}^\top \rangle, \end{aligned} \quad (35a)$$

$$\begin{aligned} \tau \dot{\mathbf{W}}_2 &= \frac{\alpha+1}{2} \beta^\top \mathbf{W}_1^\top - \mathbf{W}_2 \langle \sigma(\mathbf{W}_1 \mathbf{x}) \sigma(\mathbf{W}_1 \mathbf{x})^\top \rangle \\ &= \frac{\alpha+1}{2} \beta^\top \mathbf{W}_1^\top - \frac{\alpha+1}{2} \mathbf{W}_2 \mathbf{W}_1 \langle \mathbf{x} \sigma(\mathbf{W}_1 \mathbf{x})^\top \rangle. \end{aligned} \quad (35b)$$

We compute the second terms in the dynamics under Condition 3 and Assumption 5

$$\begin{aligned} \langle \sigma'(\mathbf{W}_1 \mathbf{x}) \odot \mathbf{W}_2^\top \mathbf{W}_2 \mathbf{W}_1 \mathbf{x} \mathbf{x}^\top \rangle &= \frac{1}{2} \langle \sigma'(\mathbf{W}_1 \mathbf{x}) \odot \mathbf{W}_2^\top \mathbf{W}_2 \mathbf{W}_1 \mathbf{x} \mathbf{x}^\top \rangle_{\mathbb{S}^+} + \frac{1}{2} \langle \sigma'(\mathbf{W}_1 \mathbf{x}) \odot \mathbf{W}_2^\top \mathbf{W}_2 \mathbf{W}_1 \mathbf{x} \mathbf{x}^\top \rangle_{\mathbb{S}^-} \\ &= \frac{1}{2} \begin{bmatrix} \mathbf{1} \\ \alpha \mathbf{1} \end{bmatrix} \odot \mathbf{W}_2^\top \mathbf{W}_2 \mathbf{W}_1 \langle \mathbf{x} \mathbf{x}^\top \rangle_{\mathbb{S}^+} + \frac{1}{2} \begin{bmatrix} \alpha \mathbf{1} \\ \mathbf{1} \end{bmatrix} \odot \mathbf{W}_2^\top \mathbf{W}_2 \mathbf{W}_1 \langle \mathbf{x} \mathbf{x}^\top \rangle_{\mathbb{S}^-} \\ &= \frac{1}{2} \begin{bmatrix} \mathbf{W}_2^{+\top} \\ \alpha \mathbf{W}_2^{-\top} \end{bmatrix} \mathbf{W}_2 \mathbf{W}_1 \Sigma + \frac{1}{2} \begin{bmatrix} \alpha \mathbf{W}_2^{+\top} \\ \mathbf{W}_2^{-\top} \end{bmatrix} \mathbf{W}_2 \mathbf{W}_1 \Sigma \\ &= \frac{\alpha+1}{2} \mathbf{W}_2^\top \mathbf{W}_2 \mathbf{W}_1 \Sigma, \end{aligned} \quad (36)$$

and

$$\begin{aligned} \langle \mathbf{x} \sigma(\mathbf{W}_1 \mathbf{x})^\top \rangle &= \frac{1}{2} \langle \mathbf{x} \sigma(\mathbf{W}_1 \mathbf{x})^\top \rangle_{\mathbb{S}^+} + \frac{1}{2} \langle \mathbf{x} \sigma(\mathbf{W}_1 \mathbf{x})^\top \rangle_{\mathbb{S}^-} \\ &= \frac{1}{2} \langle \mathbf{x} \mathbf{x}^\top \rangle_{\mathbb{S}^+} \mathbf{r} \begin{bmatrix} \mathbf{W}_2^+ & \alpha \mathbf{W}_2^- \end{bmatrix} + \frac{1}{2} \langle \mathbf{x} \mathbf{x}^\top \rangle_{\mathbb{S}^-} \mathbf{r} \begin{bmatrix} \alpha \mathbf{W}_2^+ & \mathbf{W}_2^- \end{bmatrix} \\ &= \frac{\alpha+1}{2} \Sigma \mathbf{r} \mathbf{W}_2 \\ &= \frac{\alpha+1}{2} \Sigma \mathbf{W}_1^\top. \end{aligned} \quad (37)$$

Substituting Equations (36) and (37) into Equation (35), we reduce the dynamics to

$$\begin{aligned} \tau \dot{\mathbf{W}}_1 &= \frac{\alpha+1}{2} \mathbf{W}_2^\top \beta^\top - \left(\frac{\alpha+1}{2} \right)^2 \mathbf{W}_2^\top \mathbf{W}_2 \mathbf{W}_1 \Sigma, \\ \tau \dot{\mathbf{W}}_2 &= \frac{\alpha+1}{2} \beta^\top \mathbf{W}_1^\top - \left(\frac{\alpha+1}{2} \right)^2 \mathbf{W}_2 \mathbf{W}_1 \Sigma \mathbf{W}_1^\top. \end{aligned}$$

This is the same expression as Equation (10) in the main text. If we scale the weights

$$\bar{\mathbf{W}}_1 = \sqrt{\frac{\alpha+1}{2}} \mathbf{W}_1, \quad \bar{\mathbf{W}}_2 = \sqrt{\frac{\alpha+1}{2}} \mathbf{W}_2, \quad (38)$$

the scaled (leaky) ReLU network dynamics $\bar{\mathbf{W}}(t)$ is the same as that of a linear network given in Equation (3) except for a different time constant

$$\frac{2\tau}{\alpha+1} \dot{\bar{\mathbf{W}}}_1 = \bar{\mathbf{W}}_2^\top (\beta^\top - \bar{\mathbf{W}}_2 \bar{\mathbf{W}}_1 \Sigma), \quad \frac{2\tau}{\alpha+1} \dot{\bar{\mathbf{W}}}_2 = (\beta^\top - \bar{\mathbf{W}}_2 \bar{\mathbf{W}}_1 \Sigma) \bar{\mathbf{W}}_1^\top.$$

Because Theorem 7 defines the initial condition $\bar{\mathbf{W}}(0) = \sqrt{\frac{\alpha+1}{2}} \mathbf{W}(0) = \mathbf{W}^{\text{lin}}(0)$, the weights in the linear network and the scaled weights in the (leaky) ReLU network start from the same initialization, obey the same dynamics, and consequently stay the same $\forall t \geq 0$

$$\bar{\mathbf{W}}(t) = \mathbf{W}^{\text{lin}}\left(\frac{\alpha+1}{2}t\right) \Leftrightarrow \mathbf{W}(t) = \sqrt{\frac{2}{\alpha+1}} \mathbf{W}^{\text{lin}}\left(\frac{\alpha+1}{2}t\right).$$

This proves Equation (12). Substituting Equation (12) into Equation (34) proves Equation (11)

$$\begin{aligned} f(\mathbf{x}; \mathbf{W}(t)) &= \frac{\alpha+1}{2} \mathbf{W}(t) \mathbf{W}(t) \mathbf{x} = \mathbf{W}_2^{\text{lin}}\left(\frac{\alpha+1}{2}t\right) \mathbf{W}_1^{\text{lin}}\left(\frac{\alpha+1}{2}t\right) \mathbf{x} \\ &\equiv f^{\text{lin}}\left(\mathbf{x}; \mathbf{W}^{\text{lin}}\left(\frac{\alpha+1}{2}t\right)\right). \end{aligned}$$

Part 3: We show that Assumption 5 made at time $t = 0$ remains valid for $t > 0$, meaning that weights which start with rank-one structure remain rank-one. With Assumption 5 at time $t = 0$, the dynamics of the (leaky) ReLU network is described by Equation (10). This dynamics is the same as scaled dynamics in a linear network and thus satisfies the balancing property of linear networks (Ji & Telgarsky, 2019; Du et al., 2018)

$$\frac{d}{dt} (\mathbf{W}_1 \mathbf{W}_1^\top - \mathbf{W}_2^\top \mathbf{W}_2) = 0. \quad (39)$$

Under Assumption 5 at time $t = 0$, this quantity is zero at time $t = 0$ and will stay zero

$$\forall t \geq 0: \quad \mathbf{W}_1(t) \mathbf{W}_1(t)^\top - \mathbf{W}_2(t)^\top \mathbf{W}_2(t) = \mathbf{W}_1(0) \mathbf{W}_1(0)^\top - \mathbf{W}_2(0)^\top \mathbf{W}_2(0) = \mathbf{0}.$$

Because $\text{rank}(\mathbf{W}_1 \mathbf{W}_1^\top) = \text{rank}(\mathbf{W}_1)$, the balancing property enforces that \mathbf{W}_1 and \mathbf{W}_2 have equal rank. Since \mathbf{W}_2 is a vector, \mathbf{W}_1 must also have rank one. We can write a rank-one matrix as the outer-product of two vectors $\mathbf{W}_1 = \mathbf{v} \mathbf{r}^\top$. We can assume $\|\mathbf{r}\| = 1$ for convenience and get

$$\mathbf{W}_1 \mathbf{W}_1^\top = \mathbf{v} \mathbf{r}^\top \mathbf{r} \mathbf{v}^\top = \mathbf{v} \mathbf{v}^\top = \mathbf{W}_2^\top \mathbf{W}_2 \Rightarrow \mathbf{v} = \pm \mathbf{W}_2^\top.$$

Because Assumption 5 specifies $\mathbf{W}_1 = \mathbf{W}_2^\top \mathbf{r}^\top$, then $\mathbf{v} = \mathbf{W}_2^\top$. To summarize, Assumption 5 at time $t = 0$ reduces the learning dynamics of the ReLU network to be similar to that of a linear network. The reduced dynamics satisfies the balancing property which enforces that the weights remain rank-one, thus satisfying Assumption 5 for all $t \geq 0$. \square

Proof of Theorem 7 (logistic loss). We prove Theorem 7 for logistic loss $\mathcal{L}_{\text{LG}} = \langle \ln(1 + e^{-y\hat{y}}) \rangle$.

Part 1: Same as the square loss case because proving Equation (34) does not involve the loss function.

Part 2: The gradient flow dynamics of a two-layer linear network trained with logistic loss are

$$\tau \dot{\mathbf{W}}_1^{\text{lin}} = \mathbf{W}_2^{\text{lin}\top} \mathbf{W}_2^{\text{lin}} \mathbf{W}_1^{\text{lin}} \left\langle \frac{\mathbf{x} \mathbf{x}^\top}{e^{y \mathbf{W}_2^{\text{lin}} \mathbf{W}_1^{\text{lin}} \mathbf{x}} + 1} \right\rangle, \quad (40a)$$

$$\tau \dot{\mathbf{W}}_2^{\text{lin}} = \mathbf{W}_2^{\text{lin}} \mathbf{W}_1^{\text{lin}} \left\langle \frac{\mathbf{x} \mathbf{x}^\top}{e^{y \mathbf{W}_2^{\text{lin}} \mathbf{W}_1^{\text{lin}} \mathbf{x}} + 1} \right\rangle \mathbf{W}_1^{\text{lin}\top}. \quad (40b)$$

The gradient flow dynamics of a two-layer (leaky) ReLU network trained with logistic loss are

$$\tau \dot{\mathbf{W}}_1 = \left\langle \frac{\sigma'(\mathbf{W}_1 \mathbf{x}) \odot \mathbf{W}_2^\top \mathbf{W}_2 \sigma(\mathbf{W}_1 \mathbf{x}) \mathbf{x}^\top}{e^{y \mathbf{W}_2 \sigma(\mathbf{W}_1 \mathbf{x})} + 1} \right\rangle \quad (41a)$$

$$\tau \dot{\mathbf{W}}_2 = \mathbf{W}_2 \left\langle \frac{\sigma(\mathbf{W}_1 \mathbf{x}) \sigma(\mathbf{W}_1 \mathbf{x})^\top}{e^{y \mathbf{W}_2 \sigma(\mathbf{W}_1 \mathbf{x})} + 1} \right\rangle \quad (41b)$$

Substituting Equation (34) into Equation (41), we get

$$\tau \dot{\mathbf{W}}_1 = \frac{\alpha + 1}{2} \left\langle \frac{\sigma'(\mathbf{W}_1 \mathbf{x}) \odot \mathbf{W}_2^\top \mathbf{W}_2 \mathbf{W}_1 \mathbf{x} \mathbf{x}^\top}{e^{\frac{\alpha+1}{2} y \mathbf{W}_2 \mathbf{W}_1 \mathbf{x}} + 1} \right\rangle \quad (42a)$$

$$\tau \dot{\mathbf{W}}_2 = \frac{\alpha + 1}{2} \mathbf{W}_2 \mathbf{W}_1 \left\langle \frac{\mathbf{x} \sigma(\mathbf{W}_1 \mathbf{x})^\top}{e^{\frac{\alpha+1}{2} y \mathbf{W}_2 \mathbf{W}_1 \mathbf{x}} + 1} \right\rangle \quad (42b)$$

Under Condition 3 and Assumption 5, Equation (42) can be simplified

$$\begin{aligned} & \left\langle \frac{\sigma'(\mathbf{W}_1 \mathbf{x}) \odot \mathbf{W}_2^\top \mathbf{W}_2 \mathbf{W}_1 \mathbf{x} \mathbf{x}^\top}{e^{\frac{\alpha+1}{2} y \mathbf{W}_2 \mathbf{W}_1 \mathbf{x}} + 1} \right\rangle \\ &= \frac{1}{2} \left\langle \frac{\sigma'(\mathbf{W}_1 \mathbf{x}) \odot \mathbf{W}_2^\top \mathbf{W}_2 \mathbf{W}_1 \mathbf{x} \mathbf{x}^\top}{e^{\frac{\alpha+1}{2} y \mathbf{W}_2 \mathbf{W}_1 \mathbf{x}} + 1} \right\rangle_{\mathbb{S}^+} + \frac{1}{2} \left\langle \frac{\sigma'(\mathbf{W}_1 \mathbf{x}) \odot \mathbf{W}_2^\top \mathbf{W}_2 \mathbf{W}_1 \mathbf{x} \mathbf{x}^\top}{e^{\frac{\alpha+1}{2} y \mathbf{W}_2 \mathbf{W}_1 \mathbf{x}} + 1} \right\rangle_{\mathbb{S}^-} \\ &= \frac{1}{2} \begin{bmatrix} \alpha \mathbf{W}_2^{+\top} \\ \mathbf{W}_2^{-\top} \end{bmatrix} \mathbf{W}_2 \mathbf{W}_1 \left\langle \frac{\mathbf{x} \mathbf{x}^\top}{e^{\frac{\alpha+1}{2} y \mathbf{W}_2 \mathbf{W}_1 \mathbf{x}} + 1} \right\rangle_{\mathbb{S}^+} + \frac{1}{2} \begin{bmatrix} \mathbf{W}_2^{+\top} \\ \alpha \mathbf{W}_2^{-\top} \end{bmatrix} \mathbf{W}_2 \mathbf{W}_1 \left\langle \frac{\mathbf{x} \mathbf{x}^\top}{e^{\frac{\alpha+1}{2} y \mathbf{W}_2 \mathbf{W}_1 \mathbf{x}} + 1} \right\rangle_{\mathbb{S}^-} \\ &= \frac{\alpha + 1}{2} \mathbf{W}_2^\top \mathbf{W}_2 \mathbf{W}_1 \left\langle \frac{\mathbf{x} \mathbf{x}^\top}{e^{\frac{\alpha+1}{2} y \mathbf{W}_2 \mathbf{W}_1 \mathbf{x}} + 1} \right\rangle, \end{aligned}$$

and

$$\begin{aligned} & \left\langle \frac{\mathbf{x} \sigma(\mathbf{W}_1 \mathbf{x})^\top}{e^{\frac{\alpha+1}{2} y \mathbf{W}_2 \mathbf{W}_1 \mathbf{x}} + 1} \right\rangle \\ &= \frac{1}{2} \left\langle \frac{\mathbf{x} \sigma(\mathbf{W}_1 \mathbf{x})^\top}{e^{\frac{\alpha+1}{2} y \mathbf{W}_2 \mathbf{W}_1 \mathbf{x}} + 1} \right\rangle_{\mathbb{S}^+} + \frac{1}{2} \left\langle \frac{\mathbf{x} \sigma(\mathbf{W}_1 \mathbf{x})^\top}{e^{\frac{\alpha+1}{2} y \mathbf{W}_2 \mathbf{W}_1 \mathbf{x}} + 1} \right\rangle_{\mathbb{S}^-} \\ &= \frac{1}{2} \left\langle \frac{\mathbf{x} \mathbf{x}^\top}{e^{\frac{\alpha+1}{2} y \mathbf{W}_2 \mathbf{W}_1 \mathbf{x}} + 1} \right\rangle_{\mathbb{S}^+} \mathbf{r} [\alpha \mathbf{W}_2^+ \quad \mathbf{W}_2^-] + \frac{1}{2} \left\langle \frac{\mathbf{x} \mathbf{x}^\top}{e^{\frac{\alpha+1}{2} y \mathbf{W}_2 \mathbf{W}_1 \mathbf{x}} + 1} \right\rangle_{\mathbb{S}^-} \mathbf{r} [\mathbf{W}_2^+ \quad \alpha \mathbf{W}_2^-] \\ &= \frac{\alpha + 1}{2} \left\langle \frac{\mathbf{x} \mathbf{x}^\top}{e^{\frac{\alpha+1}{2} y \mathbf{W}_2 \mathbf{W}_1 \mathbf{x}} + 1} \right\rangle \mathbf{W}_1^\top. \end{aligned}$$

Thus, the reduced dynamics of the two-layer (leaky) ReLU network are

$$\begin{aligned} \tau \dot{\mathbf{W}}_1 &= \left(\frac{\alpha + 1}{2} \right)^2 \mathbf{W}_2^\top \mathbf{W}_2 \mathbf{W}_1 \left\langle \frac{\mathbf{x} \mathbf{x}^\top}{e^{\frac{\alpha+1}{2} y \mathbf{W}_2 \mathbf{W}_1 \mathbf{x}} + 1} \right\rangle, \\ \tau \dot{\mathbf{W}}_2 &= \left(\frac{\alpha + 1}{2} \right)^2 \mathbf{W}_2 \mathbf{W}_1 \left\langle \frac{\mathbf{x} \mathbf{x}^\top}{e^{\frac{\alpha+1}{2} y \mathbf{W}_2 \mathbf{W}_1 \mathbf{x}} + 1} \right\rangle \mathbf{W}_1^\top. \end{aligned}$$

If we scale the weights as Equation (38), the (leaky) ReLU network dynamics is the same as that of a linear network given in Equation (40) except for a different time constant

$$\begin{aligned} \frac{2\tau}{\alpha + 1} \dot{\bar{\mathbf{W}}}_1 &= \bar{\mathbf{W}}_2^\top \bar{\mathbf{W}}_2 \bar{\mathbf{W}}_1 \left\langle \frac{\mathbf{x} \mathbf{x}^\top}{e^{y \bar{\mathbf{W}}_2 \bar{\mathbf{W}}_1 \mathbf{x}} + 1} \right\rangle, \\ \frac{2\tau}{\alpha + 1} \dot{\bar{\mathbf{W}}}_2 &= \bar{\mathbf{W}}_2 \bar{\mathbf{W}}_1 \left\langle \frac{\mathbf{x} \mathbf{x}^\top}{e^{y \bar{\mathbf{W}}_2 \bar{\mathbf{W}}_1 \mathbf{x}} + 1} \right\rangle \bar{\mathbf{W}}_1^\top. \end{aligned}$$

Through the same reasoning as the square loss case, Equations (11) and (12) are proved.

Part 3: Same as the square loss case. \square

C.3 Global Minimum

Proof of Corollary 9. The converged solution \mathbf{w}^* is a direct consequence of the equivalence we showed in Theorem 7 and prior results of linear networks (Saxe et al., 2014; Soudry et al., 2018).

We now show that for symmetric datasets satisfying Condition 3, the global minimum solution of a two-layer bias-free (leaky) ReLU network trained with square loss is linear.

Based on Theorem 1, we can write a two-layer bias-free (leaky) ReLU network as a linear function plus an even function $f(\mathbf{x}) = \mathbf{x}^\top \mathbf{w}^* + f_e(\mathbf{x})$ where $f_e(\cdot)$ denotes an even function. For datasets satisfying Condition 3, the square loss is

$$\begin{aligned}\mathcal{L} &= \frac{1}{2} \left\langle (y - \mathbf{x}^\top \mathbf{w}^* - f_e(\mathbf{x}))^2 \right\rangle_{p(\mathbf{x})} \\ &= \frac{1}{2} \left\langle (y - \mathbf{x}^\top \mathbf{w}^*)^2 - 2(y - \mathbf{Ax})f_e(\mathbf{x}) + f_e(\mathbf{x})^2 \right\rangle_{p(\mathbf{x})} \\ &= \frac{1}{2} \left\langle (y - \mathbf{x}^\top \mathbf{w}^*)^2 \right\rangle_{p(\mathbf{x})} + \frac{1}{2} \left\langle f_e(\mathbf{x})^2 \right\rangle_{p(\mathbf{x})}.\end{aligned}$$

The square loss attains its minimum when both $\left\langle (y - \mathbf{x}^\top \mathbf{w}^*)^2 \right\rangle_{p(\mathbf{x})}$ and $\left\langle f_e(\mathbf{x})^2 \right\rangle_{p(\mathbf{x})}$ are minimized. The former is minimized when $\mathbf{w}^* = \Sigma^{-1}\beta$. The latter is minimized when $f_e(\mathbf{x}) = 0$. Hence, for symmetric datasets satisfying Condition 3, the two-layer bias-free (leaky) ReLU network achieves globally minimum square loss with the linear, ordinary least squares solution $f(\mathbf{x}) = \mathbf{x}^\top \mathbf{w}^* = \mathbf{x}^\top \Sigma^{-1}\beta$. \square

C.4 Effect of Regularization

Theorem 7 still holds when L2 regularization is applied. Specifically, Theorem 7 holds if L2 regularization is added with hyperparameter $\lambda_\alpha = \frac{\alpha+1}{2}\lambda$, i.e., the loss is $\mathcal{L}_{\text{reg}} = \mathcal{L} + \frac{\alpha+1}{2}\lambda \|\mathbf{W}\|_2^2$. With similar calculations, we find that the regularized dynamics of the two-layer bias-free (leaky) ReLU network is

$$\tau \dot{\mathbf{W}}_1 = \frac{\alpha+1}{2} \mathbf{W}_2^\top \beta^\top - \left(\frac{\alpha+1}{2} \right)^2 \mathbf{W}_2^\top \mathbf{W}_2 \mathbf{W}_1 \Sigma - \frac{\alpha+1}{2} \lambda \mathbf{W}_1, \quad (43a)$$

$$\tau \dot{\mathbf{W}}_2 = \frac{\alpha+1}{2} \beta^\top \mathbf{W}_1^\top - \left(\frac{\alpha+1}{2} \right)^2 \mathbf{W}_2 \mathbf{W}_1 \Sigma \mathbf{W}_1^\top - \frac{\alpha+1}{2} \lambda \mathbf{W}_2. \quad (43b)$$

If we scale the weights as Equation (38), the regularized (leaky) ReLU network dynamics is again the same as that of a regularized linear network except for a different time constant

$$\frac{2\tau}{\alpha+1} \dot{\bar{\mathbf{W}}}_1 = \bar{\mathbf{W}}_2^\top (\beta^\top - \bar{\mathbf{W}}_2 \bar{\mathbf{W}}_1 \Sigma) - \lambda \bar{\mathbf{W}}_1, \quad \frac{2\tau}{\alpha+1} \dot{\bar{\mathbf{W}}}_2 = (\beta^\top - \bar{\mathbf{W}}_2 \bar{\mathbf{W}}_1 \Sigma) \bar{\mathbf{W}}_1^\top - \lambda \bar{\mathbf{W}}_2.$$

We validate this with simulations in Figure 6a. As in the unregularized case, we find that the loss curves with different leaky ReLU slopes collapse to one curve after rescaling time and the differences between weight matrices are small.

C.5 Effect of Learning Rate

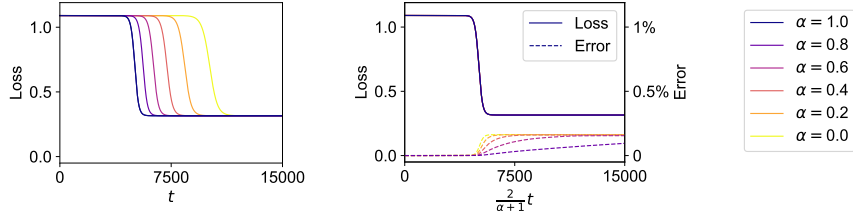
We empirically find that with a moderately large learning rate, the behaviors of two-layer bias-free (leaky) ReLU networks are consistent with Theorem 7. For simulations in Figure 6b, we use a learning rate of 0.6, which is 150 times larger than the learning rate used in Figure 2, 0.004. Due to the larger learning rate, the loss curves in Figure 6b is less smooth than those in Figures 2, 6a and 6c. Nonetheless, the loss curves with different leaky ReLU slopes collapse to one curve after rescaling time and the differences between weight matrices are small.

If the learning rate is further increased, oscillations in the loss curves occur, suggesting unstable training. In such cases, the equivalence described in Theorem 7 no longer holds. However, the learning rate is typically chosen to avoid such oscillations in training.

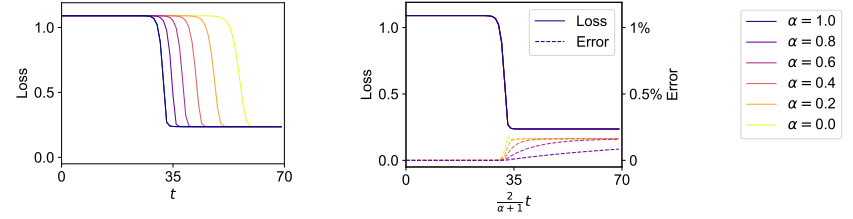
C.6 Effect of Initialization

We empirically find that under large initialization, two-layer bias-free (leaky) ReLU networks still have similar learning dynamics as its linear counterpart. As in the small initialization case, the loss curves in Figure 6c with different leaky ReLU slopes collapse to one curve after rescaling time. The differences between weight matrices are larger than those in the case of small initialization but are still less than 3%. With large initialization and a moderately large learning rate, the behavior remains consistent, as shown in Figure 6d.

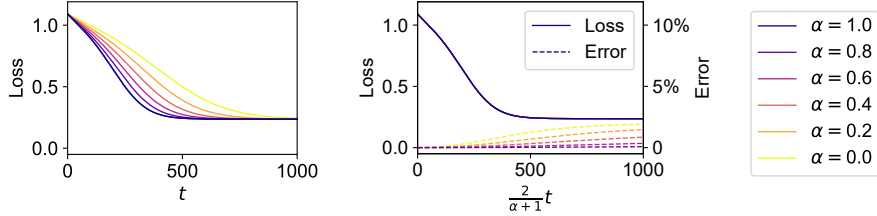
This is related to the limited expressivity of bias-free ReLU networks. Within the expressivity of two-layer bias-free ReLU networks, the linear solution is the global minimum for symmetric datasets. The two-layer bias-free ReLU network learns the linear solution starting from either small or large initialization.



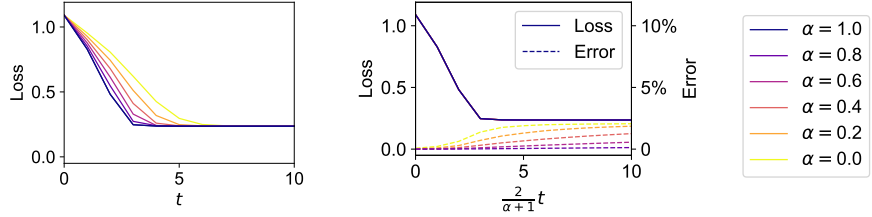
(a) Loss and error curves with L2 regularization with hyperparameter $\lambda_\alpha = 0.2(\alpha + 1)$.



(b) Loss and error curves with a moderately large learning rate, 0.6.



(c) Loss and error curves with large initialization, $w_{\text{init}} = 0.5$.



(d) Loss and error curves with large initialization, $w_{\text{init}} = 0.5$, and a moderately large learning rate, 0.6.

Figure 6: Two-layer bias-free (leaky) ReLU networks evolve like a linear network even when some of the assumptions in Theorem 7 are lifted. The setup is the same as Figure 2 except for the condition(s) specified in each individual subcaption. In (b,c,d), the time rescaling is implemented by inversely rescaling the learning rate. This avoids the inaccuracy induced by rounding the rescaled time to an integer number of epoch, which becomes non-negligible in the case of a small total number of epochs. In (c,d), with large initialization, the errors between weight matrices are larger than those in the case of small initialization but are still less than 3%.

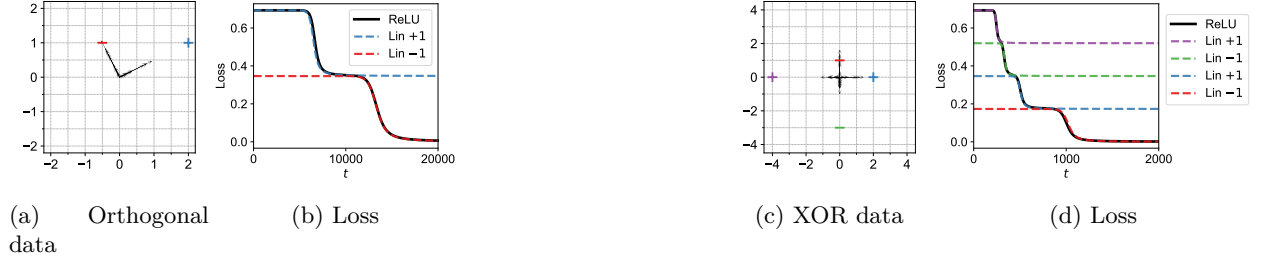


Figure 7: The same as Figure 3 but with logistic loss.

D Two-Layer Networks Learning Dynamics on Orthogonal Datasets

Suppose the low-rank weights in the two-layer bias-free ReLU network trained on the orthogonal dataset are

$$\mathbf{W}_1 = \begin{bmatrix} \mathbf{W}_1^{(1)} \\ \mathbf{W}_1^{(2)} \end{bmatrix} = \begin{bmatrix} \mathbf{r}_1^{(1)} \bar{\mathbf{x}}_1^\top \\ \mathbf{r}_1^{(2)} \bar{\mathbf{x}}_2^\top \end{bmatrix}, \quad \mathbf{W}_2 = \begin{bmatrix} \mathbf{W}_2^{(1)} & \mathbf{W}_2^{(2)} \end{bmatrix} = \begin{bmatrix} \mathbf{r}_1^{(1)\top} & \mathbf{r}_1^{(2)\top} \end{bmatrix}, \quad (44)$$

where $\bar{\mathbf{x}}_1 = \mathbf{x}_1 / \|\mathbf{x}_1\|$, $\bar{\mathbf{x}}_2 = \mathbf{x}_2 / \|\mathbf{x}_2\|$, and $\mathbf{W}_1^{(1)}, \mathbf{W}_1^{(2)}$ are rank-one block matrices and align with either one of the two data points.

Since the two data points are orthogonal $\mathbf{x}_1^\top \mathbf{x}_2 = 0$, we have for $\mu = 1, 2$

$$f(\mathbf{x}_\mu; \mathbf{W}) \equiv \mathbf{W}_2 \sigma(\mathbf{W}_1 \mathbf{x}_\mu) = \sum_{\nu=1}^2 \mathbf{W}_2^{(\nu)} \sigma(\mathbf{W}_1^{(\nu)} \mathbf{x}_\mu) = \mathbf{W}_2^{(\mu)} \sigma(\mathbf{W}_1^{(\mu)} \mathbf{x}_\mu) = \mathbf{W}_2^{(\mu)} \mathbf{W}_1^{(\mu)} \mathbf{x}_\mu$$

Learning dynamics with summed square loss for each block of weight matrices are

$$\begin{aligned} \tau \dot{\mathbf{W}}_1^{(\mu)} &= \sum_{\nu=1}^2 \sigma'(\mathbf{W}_1^{(\mu)} \mathbf{x}_\nu) \odot \mathbf{W}_2^{(\mu)\top} (y_\nu - \mathbf{W}_2 \sigma(\mathbf{W}_1 \mathbf{x}_\nu)) \mathbf{x}_\nu^\top \\ &= \sigma'(\mathbf{W}_1^{(\mu)} \mathbf{x}_\mu) \odot \mathbf{W}_2^{(\mu)\top} (y_\mu - \mathbf{W}_2 \sigma(\mathbf{W}_1 \mathbf{x}_\mu)) \mathbf{x}_\mu^\top \\ &= \mathbf{W}_2^{(\mu)\top} (y_\mu - \mathbf{W}_2^{(\mu)} \mathbf{W}_1^{(\mu)} \mathbf{x}_\mu) \mathbf{x}_\mu^\top, \\ \tau \dot{\mathbf{W}}_2^{(\mu)} &= \sum_{\nu=1}^2 (y_\nu - \mathbf{W}_2 \sigma(\mathbf{W}_1 \mathbf{x}_\nu)) \sigma(\mathbf{W}_1^{(\mu)} \mathbf{x}_\nu)^\top \\ &= (y_\mu - \mathbf{W}_2 \sigma(\mathbf{W}_1 \mathbf{x}_\mu)) \sigma(\mathbf{W}_1^{(\mu)} \mathbf{x}_\mu)^\top \\ &= (y_\mu - \mathbf{W}_2^{(\mu)} \mathbf{W}_1^{(\mu)} \mathbf{x}_\mu) \mathbf{x}_\mu^\top \mathbf{W}_1^{(\mu)\top}. \end{aligned}$$

Thus, the learning dynamics of the weights that align with \mathbf{x}_μ ($\mu = 1, 2$) are

$$\tau \dot{\mathbf{W}}_1^{(\mu)} = \mathbf{W}_2^{(\mu)\top} (y \mathbf{x}_\mu^\top - \mathbf{W}_2^{(\mu)} \mathbf{W}_1^{(\mu)} \mathbf{x}_\mu \mathbf{x}_\mu^\top), \quad (46a)$$

$$\tau \dot{\mathbf{W}}_2^{(\mu)} = (y \mathbf{x}_\mu^\top - \mathbf{W}_2^{(\mu)} \mathbf{W}_1^{(\mu)} \mathbf{x}_\mu \mathbf{x}_\mu^\top) \mathbf{W}_1^{(\mu)\top}. \quad (46b)$$

This is the same dynamics as training a linear network (with less hidden neurons) on a single data point (\mathbf{x}_μ, y_μ) as given in Equation (3).

E Deep ReLU Network Learning Dynamics

According to Equation (15), all weight elements in the intermediate layers ($\mathbf{W}_{L-1}, \dots, \mathbf{W}_3, \mathbf{W}_2$) are non-negative numbers. According to the definition of the ReLU activation function, $\sigma(\mathbf{W}_1 \mathbf{x})$ yields a vector with non-negative numbers. Thus $\mathbf{W}_2 \sigma(\mathbf{W}_1 \mathbf{x})$ yields a vector with non-negative numbers and we have $\sigma(\mathbf{W}_2 \sigma(\mathbf{W}_1 \mathbf{x})) = \mathbf{W}_2 \sigma(\mathbf{W}_1 \mathbf{x})$. Similarly, all subsequent ReLU activation functions can be ignored¹. With weights satisfying Equation (15), a deep bias-free ReLU network implements

$$f(\mathbf{x}; \mathbf{W}) \equiv \mathbf{W}_L \sigma(\dots \sigma(\mathbf{W}_2 \sigma(\mathbf{W}_1 \mathbf{x}))) = \mathbf{W}_L \dots \mathbf{W}_2 \sigma(\mathbf{W}_1 \mathbf{x}).$$

We stick to the notation for the positive and negative half-space defined in Equation (31). For $\mathbf{x} \in \mathbb{S}^+$, we have

$$f(\mathbf{x}; \mathbf{W}) = u \mathbf{W}_L \dots \mathbf{W}_2 \begin{bmatrix} \mathbf{r}_1^+ \\ \mathbf{0} \end{bmatrix} \mathbf{r}^\top \mathbf{x} = \dots = \frac{u^{L-1}}{(\sqrt{2})^{L-2}} \mathbf{W}_L \begin{bmatrix} \mathbf{r}_{L-1}^+ \\ \mathbf{0} \end{bmatrix} \mathbf{r}^\top \mathbf{x} = \left(\frac{u}{\sqrt{2}} \right)^L \mathbf{r}^\top \mathbf{x}.$$

For $\mathbf{x} \in \mathbb{S}^-$, we have

$$f(\mathbf{x}; \mathbf{W}) = u \mathbf{W}_L \dots \mathbf{W}_2 \begin{bmatrix} \mathbf{0} \\ \mathbf{r}_1^- \end{bmatrix} \mathbf{r}^\top \mathbf{x} = \dots = \frac{u^{L-1}}{(\sqrt{2})^{L-2}} \mathbf{W}_L \begin{bmatrix} \mathbf{0} \\ \mathbf{r}_{L-1}^- \end{bmatrix} \mathbf{r}^\top \mathbf{x} = \left(\frac{u}{\sqrt{2}} \right)^L \mathbf{r}^\top \mathbf{x}.$$

Hence, the deep bias-free ReLU network implements a linear function $f(\mathbf{x}; \mathbf{W}) = \left(\frac{u}{\sqrt{2}} \right)^L \mathbf{r}^\top \mathbf{x}$. Notice that a deep linear network with such weights implement

$$\mathbf{W}_L \dots \mathbf{W}_2 \mathbf{W}_1 \mathbf{x} = u \mathbf{W}_L \dots \mathbf{W}_2 \begin{bmatrix} \mathbf{r}_1^+ \\ \mathbf{r}_1^- \end{bmatrix} \mathbf{r}^\top \mathbf{x} = \dots = \frac{u^{L-1}}{(\sqrt{2})^{L-2}} \mathbf{W}_L \mathbf{r}_{L-1} \mathbf{r}^\top \mathbf{x} = \frac{u^L}{(\sqrt{2})^{L-2}} \mathbf{r}^\top \mathbf{x}.$$

Equation (16) is thus proved.

We now prove that under Equation (15) on the weights at time $t = 0$, we have that $\forall t \geq 0$, Equation (15) remains valid. We assume that $\sigma'(0) = 0$. We substitute the low-rank weights defined in Equation (15) into the learning dynamics of deep bias-free ReLU networks and make simplifications. For the first layer,

$$\begin{aligned} \tau \dot{\mathbf{W}}_1 &= \langle \sigma'(\mathbf{W}_1 \mathbf{x}) \odot \mathbf{W}_2^\top \dots \mathbf{W}_L^\top (y - f(\mathbf{x})) \mathbf{x}^\top \rangle \\ &= \frac{u^{L-1}}{(\sqrt{2})^{L-2}} \left\langle \sigma'(\mathbf{W}_1 \mathbf{x}) \odot \mathbf{r}_1 \left(y - \left(\frac{u}{\sqrt{2}} \right)^L \mathbf{r}^\top \mathbf{x} \right) \mathbf{x}^\top \right\rangle \\ &= \frac{u^{L-1}}{(\sqrt{2})^L} \begin{bmatrix} \mathbf{r}_1^+ \\ \mathbf{0} \end{bmatrix} \left\langle \left(y - \left(\frac{u}{\sqrt{2}} \right)^L \mathbf{r}^\top \mathbf{x} \right) \mathbf{x}^\top \right\rangle_{\mathbb{S}^+} + \frac{u^{L-1}}{(\sqrt{2})^L} \begin{bmatrix} \mathbf{0} \\ \mathbf{r}_1^- \end{bmatrix} \left\langle \left(y - \left(\frac{u}{\sqrt{2}} \right)^L \mathbf{r}^\top \mathbf{x} \right) \mathbf{x}^\top \right\rangle_{\mathbb{S}^-} \\ &= \frac{u^{L-1}}{(\sqrt{2})^L} \mathbf{r}_1 \left(\beta^\top - \left(\frac{u}{\sqrt{2}} \right)^L \mathbf{r}^\top \Sigma \right). \end{aligned} \quad (47)$$

¹The activation functions can be ignored when calculating the network output but cannot be ignored when calculating the gradients. This is because for a ReLU function $\sigma(z) = \max(z, 0)$ and a linear function $\phi(z) = z$, the function values are equal at zero $\sigma(0) = \phi(0)$ but the derivatives are not equal at zero $\sigma'(0) \neq \phi'(0)$.

For intermediate layers $1 < l < L$,

$$\begin{aligned}
\tau \dot{\mathbf{W}}_l &= \langle \sigma'(\mathbf{h}_l) \odot \mathbf{W}_{l+1}^\top \cdots \mathbf{W}_L^\top (y - f(\mathbf{x})) \sigma(\mathbf{h}_{l-1})^\top \rangle \\
&= \left(\frac{u}{\sqrt{2}} \right)^{L-1} \begin{bmatrix} \mathbf{1} \\ \mathbf{0} \end{bmatrix} \odot \begin{bmatrix} \mathbf{r}_l^+ \\ \mathbf{r}_l^- \end{bmatrix} \left\langle \left(y - \left(\frac{u}{\sqrt{2}} \right)^L \mathbf{r}^\top \mathbf{x} \right) \mathbf{x}^\top \right\rangle_{\mathbb{S}^+} \mathbf{r} \begin{bmatrix} \mathbf{r}_{l-1}^{+\top} & \mathbf{0} \end{bmatrix} \\
&\quad + \left(\frac{u}{\sqrt{2}} \right)^{L-1} \begin{bmatrix} \mathbf{0} \\ \mathbf{1} \end{bmatrix} \odot \begin{bmatrix} \mathbf{r}_l^+ \\ \mathbf{r}_l^- \end{bmatrix} \left\langle \left(y - \left(\frac{u}{\sqrt{2}} \right)^L \mathbf{r}^\top \mathbf{x} \right) \mathbf{x}^\top \right\rangle_{\mathbb{S}^-} \mathbf{r} \begin{bmatrix} \mathbf{0} & \mathbf{r}_{l-1}^{-\top} \end{bmatrix} \\
&= \left(\frac{u}{\sqrt{2}} \right)^{L-1} \begin{bmatrix} \mathbf{r}_l^+ \mathbf{r}_{l-1}^{+\top} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \left(\beta^\top - \left(\frac{u}{\sqrt{2}} \right)^L \mathbf{r}^\top \Sigma \right) \mathbf{r} \\
&\quad + \left(\frac{u}{\sqrt{2}} \right)^{L-1} \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{r}_l^- \mathbf{r}_{l-1}^{-\top} \end{bmatrix} \left(\beta^\top - \left(\frac{u}{\sqrt{2}} \right)^L \mathbf{r}^\top \Sigma \right) \mathbf{r} \\
&= \frac{u^{L-1}}{(\sqrt{2})^L} \begin{bmatrix} \sqrt{2} \mathbf{r}_l^+ \mathbf{r}_{l-1}^{+\top} & \mathbf{0} \\ \mathbf{0} & \sqrt{2} \mathbf{r}_l^- \mathbf{r}_{l-1}^{-\top} \end{bmatrix} \left(\beta^\top - \left(\frac{u}{\sqrt{2}} \right)^L \mathbf{r}^\top \Sigma \right) \mathbf{r}. \tag{48}
\end{aligned}$$

For the last layer,

$$\begin{aligned}
\tau \dot{\mathbf{W}}_L &= \langle (y - f(\mathbf{x})) \sigma(\mathbf{h}_{L-1})^\top \rangle \\
&= \frac{u^{L-1}}{(\sqrt{2})^L} \left\langle \left(y - \left(\frac{u}{\sqrt{2}} \right)^L \mathbf{r}^\top \mathbf{x} \right) \mathbf{x}^\top \right\rangle_{\mathbb{S}^+} \mathbf{r} \begin{bmatrix} \mathbf{r}_{L-1}^{+\top} & \mathbf{0} \end{bmatrix} \\
&\quad + \frac{u^{L-1}}{(\sqrt{2})^L} \left\langle \left(y - \left(\frac{u}{\sqrt{2}} \right)^L \mathbf{r}^\top \mathbf{x} \right) \mathbf{x}^\top \right\rangle_{\mathbb{S}^-} \mathbf{r} \begin{bmatrix} \mathbf{0} & \mathbf{r}_{L-1}^{-\top} \end{bmatrix} \\
&= \frac{u^{L-1}}{(\sqrt{2})^L} \left(\beta^\top - \left(\frac{u}{\sqrt{2}} \right)^L \mathbf{r}^\top \Sigma \right) \mathbf{r} \mathbf{r}_{L-1}^\top. \tag{49}
\end{aligned}$$

Equations (47) to (49) can be written more compactly as Equation (17) if we substitute the weights back in. The dynamics of the deep ReLU network as in Equation (17) is the same as a deep linear network as in Equation (7) except for constant coefficients.

We now prove that the low-rank weights remain low-rank once formed. We substitute the low-rank weights defined in Equation (15) into the left-hand side of Equations (47) to (49) and get

$$\begin{aligned}
\tau \frac{d}{dt} u \mathbf{r}_1 \mathbf{r}^\top &= \frac{u^{L-1}}{(\sqrt{2})^L} \mathbf{r}_1 \left(\beta^\top - \left(\frac{u}{\sqrt{2}} \right)^L \mathbf{r}^\top \Sigma \right), \\
\tau \frac{d}{dt} u \begin{bmatrix} \sqrt{2} \mathbf{r}_l^+ \mathbf{r}_{l-1}^{+\top} & \mathbf{0} \\ \mathbf{0} & \sqrt{2} \mathbf{r}_l^- \mathbf{r}_{l-1}^{-\top} \end{bmatrix} &= \frac{u^{L-1}}{(\sqrt{2})^L} \begin{bmatrix} \sqrt{2} \mathbf{r}_l^+ \mathbf{r}_{l-1}^{+\top} & \mathbf{0} \\ \mathbf{0} & \sqrt{2} \mathbf{r}_l^- \mathbf{r}_{l-1}^{-\top} \end{bmatrix} \left(\beta^\top - \left(\frac{u}{\sqrt{2}} \right)^L \mathbf{r}^\top \Sigma \right) \mathbf{r}, \\
\tau \frac{d}{dt} u \mathbf{r}_{L-1}^\top &= \frac{u^{L-1}}{(\sqrt{2})^L} \left(\beta^\top - \left(\frac{u}{\sqrt{2}} \right)^L \mathbf{r}^\top \Sigma \right) \mathbf{r} \mathbf{r}_{L-1}^\top.
\end{aligned}$$

We cancel out the nonzero common terms on both sides and reduce the dynamics to two differential equations. The first one is about the norm of a layer u . The second one is about the rank-one direction in the first layer $u \mathbf{r}$.

$$\begin{aligned}
\tau \frac{d}{dt} u &= \frac{u^{L-1}}{(\sqrt{2})^L} \left(\beta^\top - \left(\frac{u}{\sqrt{2}} \right)^L \mathbf{r}^\top \Sigma \right) \mathbf{r}, \\
\tau \frac{d}{dt} u \mathbf{r}^\top &= \frac{u^{L-1}}{(\sqrt{2})^L} \left(\beta^\top - \left(\frac{u}{\sqrt{2}} \right)^L \mathbf{r}^\top \Sigma \right).
\end{aligned}$$

After the weights have formed the low-rank structure specified in Equation (15), the norm of each layer u and the rank-one direction of the first layer \mathbf{r} evolve while $\mathbf{r}_1, \mathbf{r}_2, \dots, \mathbf{r}_{L-1}$ stay fixed. Hence, under Equation (15) on the weights at time $t = 0$, we have that $\forall t \geq 0$, Equation (15) remains valid.

F Depth Separation

We visualize $g(\mathbf{x})$ defined in Equation (9) below. This function cannot be expressed by any two-layer bias-free ReLU network because it is odd and nonlinear. Corollary 2 indicates that odd and nonlinear functions cannot be expressed by two-layer bias-free ReLU networks. It can be expressed by a three-layer bias-free ReLU network as written in Equation (9).

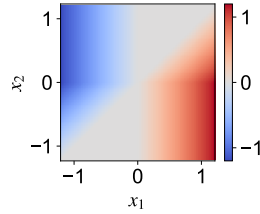


Figure 8: Function $g(\mathbf{x})$ defined in Equation (9) is plotted with color.

A clarification on Section 3.2 is that deep bias-free ReLU networks are not more expressive than their two-layer counterparts if the input is scalar. For scalar input functions, the only positively homogeneous odd function is the linear function. Neither two-layer nor deep bias-free ReLU networks can express nonlinear odd functions with scalar input.

Another relevant fact is that there is also depth separation between two-layer and deep ReLU networks with bias. One example is the pyramid function, $\sigma(1 - \|\mathbf{x}\|_1)$, which was studied in (Ongie et al., 2020, Example 4) and (Nacson et al., 2023, Proposition 2)

G Additional Figure

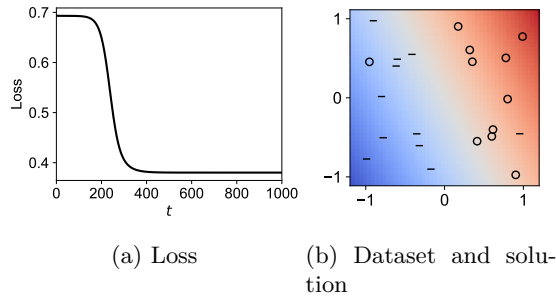


Figure 9: Two-layer bias-free ReLU network trained on a linearly separable binary classification task with label flipping noise. Since the dataset satisfies symmetric Condition 3, the network follows linear network dynamics and converges to a linear decision boundary, which is a presumably robust solution here as it avoids overfitting the two noisy labels. (a) Loss curve. Logistic loss is used here. (b) The dataset is plotted with empty circles and short lines, representing data points with $+1$ labels and -1 labels respectively. The network output at the end of training is plotted in color.

H Implementation Details

All networks are initialized with small random weights. Specifically, the initial weights in the l -th layer are sampled i.i.d. from a normal distribution $\mathcal{N}(0, w_{\text{init}}^2/N_l)$ where N_l is the number of weight parameters in the l -th layer. The initialization scale w_{init} is specified below.

Figure 1. The networks have width 100. The initialization scale $w_{\text{init}} = 10^{-2}$. The learning rate is 0.2. The two-layer networks are trained 10000 epochs. The three-layer networks are trained 80000 epochs. The dataset is plotted in the figure. Empty black circles denote data points with +1 labels. Short black lines denote data points with -1 labels. The size of the datasets is 120.

Figure 2. The networks have width 500. The initialization scale is $w_{\text{init}} = 10^{-8}$. The learning rate is 0.004. The input is 20-dimensional, $\mathbf{x} \in \mathbb{R}^{20}$. We sample 1000 i.i.d. vectors $\mathbf{x}_n \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ and include both \mathbf{x}_n and $-\mathbf{x}_n$ in the dataset, resulting in 2000 data points. The output is generated as $y = \mathbf{w}^\top \mathbf{x} + \sin(4\mathbf{w}^\top \mathbf{x})$ where elements of \mathbf{w} are randomly sampled from a uniform distribution $\mathcal{U}[-0.5, 0.5]$. This dataset satisfy Condition 3 since the empirical input distribution is even and the output is generated by an odd function.

Figures 3 and 7. We use the same hyperparameters as Boursier et al. (2022). The network width is 60. The initialization scale $w_{\text{init}} = 10^{-6}$. The learning rate is 0.001 for square loss and 0.004 for logistic loss. The orthogonal input dataset contains two data points, i.e., $[-0.5, 1]$, $[2, 1]$. The XOR input dataset contains four data points, i.e., $[0, 1]$, $[2, 0]$, $[0, -3]$, $[-4, 0]$.

Figure 4. The networks have width 100. The initialization scale $w_{\text{init}} = 10^{-2}$. The learning rate is 0.1. The networks are trained 20000 epochs. The dataset is generated in the same way as Figure 2 except that the output is generated as $y = \mathbf{w}^\top \mathbf{x}$.

Figure 5. The network width is 100. The initialization scale $w_{\text{init}} = 10^{-3}$. The learning rate is 0.025. The dataset contains six data points: $[1, 1]$, $[-1, -1]$, $[1, -1]$, $[-1, 1]$, $[-1, 0]$, $[1, \Delta y]$.