

UNDERSTANDING OVERPARAMETERIZATION IN GENERATIVE ADVERSARIAL NETWORKS

Yogesh Balaji^{1*}, Mohammadmahdi Sajedi^{2*}, Neha Mukund Kalibhat¹, Mucong Ding¹,
Dominik Stöger², Mahdi Soltanolkotabi², Soheil Feizi¹

¹ University of Maryland, College Park, MD

² University of Southern California, Los Angeles, CA

ABSTRACT

A broad class of *unsupervised* deep learning methods such as Generative Adversarial Networks (GANs) involve training of overparameterized models where the number of parameters of the model exceeds a certain threshold. Indeed, most successful GANs used in practice are trained using overparameterized generator and discriminator networks, both in terms of depth and width. A large body of work in *supervised* learning have shown the importance of model overparameterization in the convergence of the gradient descent (GD) to globally optimal solutions. In contrast, the unsupervised setting and GANs in particular involve non-convex concave mini-max optimization problems that are often trained using Gradient Descent/Ascent (GDA). The role and benefits of model overparameterization in the convergence of GDA to a global saddle point in non-convex concave problems is far less understood. In this work, we present a comprehensive analysis of the importance of model overparameterization in GANs both theoretically and empirically. We theoretically show that in an overparameterized GAN model with a 1-layer neural network generator and a linear discriminator, GDA converges to a global saddle point of the underlying non-convex concave min-max problem. To the best of our knowledge, this is the first result for global convergence of GDA in such settings. Our theory is based on a more general result that holds for a broader class of nonlinear generators and discriminators that obey certain assumptions (including deeper generators and random feature discriminators). Our theory utilizes and builds upon a novel connection with the convergence analysis of linear time-varying dynamical systems which may have broader implications for understanding the convergence behavior of GDA for non-convex concave problems involving overparameterized models. We also empirically study the role of model overparameterization in GANs using several large-scale experiments on CIFAR-10 and Celeb-A datasets. Our experiments show that overparameterization improves the quality of generated samples across various model architectures and datasets. Remarkably, we observe that overparameterization leads to faster and more stable convergence behavior of GDA across the board.

1 INTRODUCTION

In recent years, we have witnessed tremendous progress in deep generative modeling with some state-of-the-art models capable of generating photo-realistic images of objects and scenes (Brock et al., 2019; Karras et al., 2019; Clark et al., 2019). Three prominent classes of deep generative models include GANs (Goodfellow et al., 2014), VAEs (Kingma & Welling, 2014) and normalizing flows (Dinh et al., 2017). Of these, GANs remain a popular choice for data synthesis especially in the image domain. GANs are based on a two player *min-max* game between a generator network that generates samples from a distribution, and a critic (discriminator) network that discriminates real distribution from the generated one. The networks are optimized using Gradient Descent/Ascent (GDA) to reach a saddle-point of the min-max optimization problem.

*First two authors contributed equally. Correspondence to yogesh@cs.umd.edu, sajedi@usc.edu

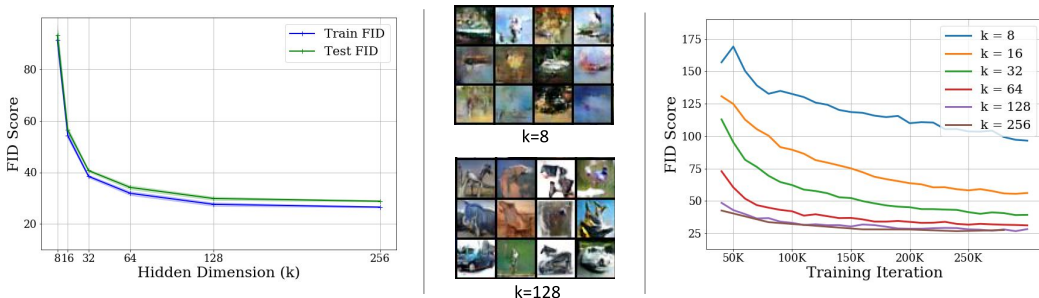


Figure 1: **Overparameterization in GANs.** We train DCGAN models by varying the size of the hidden dimension k (larger the k , more overparameterized the models are, see Fig. 8 for details). Overparameterized GANs enjoy improved training and test FID scores (*the left panel*), generate high-quality samples (*the middle panel*) and have fast and stable convergence (*the right panel*).

One of the key factors that has contributed to the successful training of GANs is model *overparameterization*, defined based on the model parameters count. By increasing the complexity of discriminator and generator networks, both in depth and width, recent papers show that GANs can achieve photo-realistic image and video synthesis (Brock et al., 2019; Clark et al., 2019; Karras et al., 2019). While these works empirically demonstrate some benefits of *overparameterization*, there is lack of a rigorous study explaining this phenomena. In this work, we attempt to provide a comprehensive understanding of the role of *overparameterization* in GANs, both theoretically and empirically. We note that while *overparameterization* is a key factor in training successful GANs, other factors such as generator and discriminator architectures, regularization functions and model hyperparameters have to be taken into account as well to improve the performance of GANs.

Recently, there has been a large body of work in *supervised* learning (e.g. regression or classification problems) studying the importance of model overparameterization in gradient descent (GD)’s convergence to globally optimal solutions (Soltanolkotabi et al., 2018; Allen-Zhu et al., 2019; Du et al., 2019; Oymak & Soltanolkotabi, 2019; Zou & Gu, 2019; Oymak et al., 2019). A key observation in these works is that, under some conditions, overparameterized models experience *lazy training* (Chizat et al., 2019) where optimal model parameters computed by GD remain close to a randomly initialized model. Thus, using a linear approximation of the model in the parameter space, one can show the global convergence of GD in such minimization problems.

In contrast, training GANs often involves solving a non-convex concave *min-max* optimization problem that fundamentally differs from a single minimization problem of classification/regression. The key question is whether overparameterized GANs also experience lazy training in the sense that overparameterized generator and discriminator networks remain sufficiently close to their initializations. This may then lead to a general theory of global convergence of GDA for such overparameterized non-convex concave min-max problems.

In this paper we first theoretically study the role of overparameterization for a GAN model with a 1-hidden layer generator and a linear discriminator. We study two optimization procedures to solve this problem: (i) using a conventional training procedure in GANs based on GDA in which generator and discriminator networks perform simultaneous steps of gradient descent to optimize their respective models, (ii) using GD to optimize generator’s parameters for the optimal discriminator. The latter case corresponds to taking a sufficiently large number of gradient ascent steps to update discriminator’s parameters for each GD step of the generator. In both cases, our results show that in an overparameterized regime, the GAN optimization converges to a global solution. To the best of our knowledge, this is the first result showing the global convergence of GDA in such settings. While in our results we focus on one-hidden layer generators and linear discriminators, our theory is based on analyzing a general class of min-max optimization problems which can be used to study a much broader class of generators and discriminators potentially including deep generators and deep random feature-based discriminators. A key component of our analysis is a novel connection to exponential stability of non-symmetric time varying dynamical systems in control theory which may have broader implications for theoretical analysis of GAN’s training. Ideas from control theory have

also been used for understanding and improving training dynamics of GANs in (Xu et al., 2019; An et al., 2018).

Having analyzed overparameterized GANs for relatively simple models, we next provide a comprehensive empirical study of this problem for practical GANs such as DCGAN (Radford et al., 2016) and ResNet GAN (Gulrajani et al., 2017) trained on CIFAR-10 and Celeb-A datasets. For example, the benefit of overparameterization in training DCGANs on CIFAR-10 is illustrated in Figure 1. We have three key observations: (i) as the model becomes more overparameterized (e.g. using wider networks), the *training* FID scores that measure the training error, decrease. This phenomenon has been observed in other studies as well (Brock et al., 2019). (ii) overparameterization does not hurt the *test* FID scores (i.e. the generalization gap remains small). This improved test-time performance can also be seen qualitatively in the center panel of Figure 1, where overparameterized models produce samples of improved quality. (iii) Remarkably, overparameterized GANs, with a lot of parameters to optimize over, have significantly improved convergence behavior of GDA, both in terms of rate and stability, compared to small GAN models (see the right panel of Figure 1).

In summary, in this paper

- We provide the first theoretical guarantee of simultaneous GDA’s global convergence for an overparameterized GAN with one-hidden neural network generator and a linear discriminator (Theorem 2.1).
- By establishing connections with linear time-varying dynamical systems, we provide a theoretical framework to analyze simultaneous GDA’s global convergence for a general overparameterized GAN (including deeper generators and random feature discriminators), under some general conditions (Theorems 2.3 and A.4).
- We provide a comprehensive empirical study of the role of model overparameterization in GANs using several large-scale experiments on CIFAR-10 and Celeb-A datasets. We observe overparameterization improves GANs’ training error, generalization error, sample qualities as well as the convergence rate and stability of GDA.

2 THEORETICAL RESULTS

2.1 PROBLEM FORMULATION

Given n data points of the form $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n \in \mathbb{R}^m$, the goal of GAN’s training is to find a generator that can mimic sampling from the same distribution as the training data. More specifically, the goal is to find a generator mapping $\mathcal{G}_\theta(z) : \mathbb{R}^d \rightarrow \mathbb{R}^m$, parameterized by $\theta \in \mathbb{R}^p$, so that $\mathcal{G}_\theta(z_1), \mathcal{G}_\theta(z_2), \dots, \mathcal{G}_\theta(z_n)$ with z_1, z_2, \dots, z_n generated i.i.d. according to $\mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ has a similar empirical distribution to $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ ¹. To measure the discrepancy between the data points and the GAN outputs, one typically uses a discriminator mapping $\mathcal{D}_{\tilde{\theta}} : \mathbb{R}^m \rightarrow \mathbb{R}$ parameterized with $\tilde{\theta} \in \mathbb{R}^{\tilde{p}}$. The overall training approach takes the form of the following min-max optimization problem which minimizes the worst-case discrepancy detected by the discriminator

$$\min_{\theta} \max_{\tilde{\theta}} \frac{1}{n} \sum_{i=1}^n \mathcal{D}_{\tilde{\theta}}(\mathbf{x}_i) - \frac{1}{n} \sum_{i=1}^n \mathcal{D}_{\tilde{\theta}}(\mathcal{G}_\theta(z_i)) + \mathcal{R}(\tilde{\theta}). \quad (1)$$

Here, $\mathcal{R}(\tilde{\theta})$ is a regularizer that typically ensures the discriminator is Lipschitz. This formulation mimics the popular Wasserstein GAN (Arjovsky et al., 2017) (or, IPM GAN) formulations. This optimization problem is typically solved by running Gradient Descent Ascent (GDA) on the minimization/maximization variables.

The generator and discriminator mappings \mathcal{G} and \mathcal{D} used in practice are often deep neural networks. Thus, the min-max optimization problem above is highly nonlinear and non-convex concave. Saddle point optimization is a classical and fundamental problem in game theory (Von Neumann & Morgenstern, 1953) and control (Gutman, 1979). However, most of the classical results apply to the

¹In general, the number of observed and generated samples can be different. However, in practical GAN implementations, batch sizes of observed and generated samples are usually the same. Thus, for simplicity, we make this assumption in our setup.

convex-concave case (Arrow et al., 1958) while the saddle point optimization of GANs is often *non convex-concave*. If GDA converges to the global (local) saddle points, we say it is globally (locally) stable. For a general min-max optimization, however, GDA can be trapped in a loop or even diverge. Except in some special cases (e.g. (Feizi et al., 2018) for a quadratic GAN formulation or (Lei et al., 2019) for the under-parametrized setup when the generator is a one-layer network), GDA is not globally stable for GANs in general (Nagarajan & Kolter, 2017; Mescheder et al., 2018; Adolphs et al., 2019; Mescheder et al., 2017; Daskalakis et al., 2020).

None of these works, however, study the role of model overparameterization in the global/local convergence (stability) of GDA. In particular, it has been empirically observed (as we also demonstrate in this paper) that when the generator/discriminator contain a large number of parameters (i.e. are sufficiently overparameterized) GDA does indeed find (near) globally optimal solutions. In this section we wish to demystify this phenomenon from a theoretical perspective.

2.2 DEFINITION OF MODEL OVERPARAMETERIZATION

In this paper, we use *overparameterization* in the context of model parameters count. Informally speaking, overparameterized models have large number of parameters, that is we assume that the number of model parameters is sufficiently large. In specific problem setups of Section 2, we precisely compute thresholds where the number of model parameters should exceed in order to observe nice convergence properties of GDA. Note that the definition of *overparameterization* based on model parameters count is related, but distinct from the complexity of the hypothesis class. For instance, in our empirical studies, when we say we *overparameterize* a neural network, we fix the number of layers in the neural network and increase the hidden dimensions. Our definition does not include the case where the number of layers also increases, which forms a different hypothesis class.

2.3 RESULTS FOR ONE-HIDDEN LAYER GENERATORS AND RANDOM DISCRIMINATORS

In this section, we discuss our main results on the convergence of gradient based algorithms when training GANs in the overparameterized regime. We focus on the case where the generator takes the form of a single hidden-layer ReLU network with d inputs, k hidden units, and m outputs. Specifically, $\mathcal{G}(\mathbf{z}) = \mathbf{V} \cdot \text{ReLU}(\mathbf{W}\mathbf{z})$ with $\mathbf{W} \in \mathbb{R}^{k \times d}$ and $\mathbf{V} \in \mathbb{R}^{m \times k}$ denoting the input-to-hidden and hidden-to-output weights. We also consider a linear discriminator of the form $\mathcal{D}(\mathbf{x}) = \mathbf{d}^T \mathbf{x}$ with an ℓ_2 regularizer on the weights i.e. $\mathcal{R}(\mathbf{d}) = -\|\mathbf{d}\|_{\ell_2}^2/2$. The overall min-max optimization problem (equation 1) takes the form

$$\min_{\mathbf{W} \in \mathbb{R}^{k \times d}} \max_{\mathbf{d} \in \mathbb{R}^m} \mathcal{L}(\mathbf{W}, \mathbf{d}) := \langle \mathbf{d}, \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \mathbf{V} \text{ReLU}(\mathbf{W}\mathbf{z}_i)) \rangle - \frac{\|\mathbf{d}\|_{\ell_2}^2}{2}. \quad (2)$$

Note that we initialize \mathbf{V} at random and keep it fixed throughout the training. The common approach to solve the above optimization problem is to run a Gradient Descent Ascent (GDA) algorithm. At iteration t , GDA takes the form

$$\begin{cases} \mathbf{d}_{t+1} &= \mathbf{d}_t + \mu \nabla_{\mathbf{d}} \mathcal{L}(\mathbf{W}_t, \mathbf{d}_t) \\ \mathbf{W}_{t+1} &= \mathbf{W}_t - \eta \nabla_{\mathbf{W}} \mathcal{L}(\mathbf{W}_t, \mathbf{d}_t) \end{cases} \quad (3)$$

Next, we establish the global convergence of GDA for an overparameterized model. Note that a global saddle point $(\mathbf{W}^*, \mathbf{d}^*)$ is defined as

$$\mathcal{L}(\mathbf{W}^*, \mathbf{d}) \leq \mathcal{L}(\mathbf{W}^*, \mathbf{d}^*) \leq \mathcal{L}(\mathbf{W}, \mathbf{d}^*)$$

for all feasible \mathbf{W} and \mathbf{d} . If these inequalities hold in a local neighborhood, $(\mathbf{W}^*, \mathbf{d}^*)$ is called a local saddle point.

Theorem 2.1 *Let $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n \in \mathbb{R}^m$ be n training data with their mean defined as $\bar{\mathbf{x}} := \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i$. Consider the GAN model with a linear discriminator of the form $\mathcal{D}(\mathbf{x}) = \mathbf{d}^T \mathbf{x}$ parameterized by $\mathbf{d} \in \mathbb{R}^m$ and a one hidden layer neural network generator of the form $\mathcal{G}(\mathbf{z}) = \mathbf{V} \phi(\mathbf{W}\mathbf{z})$ parameterized by $\mathbf{W} \in \mathbb{R}^{k \times d}$ with $\mathbf{V} \in \mathbb{R}^{m \times k}$ a fixed matrix generated at random with i.i.d. $\mathcal{N}(0, \sigma_v^2)$ entries. Also assume the input data to the generator $\{\mathbf{z}_i\}_{i=1}^n$ are generated i.i.d. according to $\sim \mathcal{N}(\mathbf{0}, \sigma_z^2 \mathbf{I}_d)$. Furthermore, assume the generator weights at initialization $\mathbf{W}_0 \in \mathbb{R}^{k \times d}$*

are generated i.i.d. according to $\mathcal{N}(0, \sigma_w^2)$. Furthermore, assume the standard deviations above obey $\sigma_v \sigma_w \sigma_z \geq \|\bar{\mathbf{x}}\|_{\ell_2} / (md^{\frac{5}{2}} \log d^{\frac{3}{2}})$. Then, as long as

$$k \geq C \cdot md^4 \log(d)^3$$

with C a fixed constant, running GDA updates per equation 3 starting from the random \mathbf{W}_0 above and $\mathbf{d}_0 = \mathbf{0}^2$ with step-sizes obeying $0 < \mu \leq 1$ and $\eta = \bar{\eta} \frac{\frac{\mu}{n-1}}{324 \cdot k \cdot \frac{d+\frac{n-1}{n}}{n} \cdot \sigma_v^2 \cdot \sigma_z^2}$, with $\bar{\eta} \leq 1$, satisfies

$$\left\| \frac{1}{n} \sum_{i=1}^n \text{VReLU}(\mathbf{W}_\tau \mathbf{z}_i) - \bar{\mathbf{x}} \right\|_{\ell_2} \leq 5 (1 - 10^{-5} \cdot \bar{\eta} \mu)^\tau \left\| \frac{1}{n} \sum_{i=1}^n \text{VReLU}(\mathbf{W}_0 \mathbf{z}_i) - \bar{\mathbf{x}} \right\|_{\ell_2}. \quad (4)$$

This holds with probability at least $1 - (n+5)e^{-\frac{m}{1500}} - 5k \cdot e^{-c_1 \cdot n} - (2k+2)e^{-\frac{d}{216}} - ne^{-c_2 \cdot md^3 \log(d)^2}$ where c_1, c_2 are fixed numerical constants.

To better understand the implications of the above theorem, note that the objective of equation 2 can be simplified by solving the inner maximization in a closed form so that the min-max problem in equation 2 is equivalent to the following single minimization problem:

$$\min_{\mathbf{W}} \mathcal{L}(\mathbf{W}) := \frac{1}{2} \left\| \frac{1}{n} \sum_{i=1}^n \text{VReLU}(\mathbf{W} \mathbf{z}_i) - \bar{\mathbf{x}} \right\|_{\ell_2}^2, \quad (5)$$

which has a global optimum of zero. As a result equation 4 in Theorem 2.1 guarantees that running simultaneous GDA updates achieves the global optimum. This holds as long as the generator network is sufficiently overparameterized in the sense that the number of hidden nodes is polynomially large in its output dimension m and input dimension d . Interestingly, the rate of convergence guaranteed by this result is geometric, guaranteeing fast GDA convergence to the global optima. To the extent of our knowledge, this is the first result that establishes the global convergence of simultaneous GDA for an overparameterized GAN model.

While the result proved above shows the global convergence of GDA for a GAN with 1-hidden layer generator and a linear discriminator, for a general GAN model, local saddle points may not even exist and GDA may converge to approximate local saddle points (Berard et al., 2020; Farnia & Ozdaglar, 2020). For a general min-max problem, (Daskalakis et al., 2020) has recently shown that *approximate* local saddle points exist under some general conditions on the lipschitzness of the objective function. Understanding GDA dynamics for a general GAN remains an important open problem. Our result in Theorem 2.1 is a first and important step towards that.

We acknowledge that the considered GAN formulation of equation 2 is very simpler than GANs used in practice. Specially, since the discriminator is linear, this GAN can be viewed as a moment-matching GAN (Li et al., 2017) pushing first moments of input and generative distributions towards each other. Alternatively, this GAN formulation can be viewed as one instance of the Sliced Wasserstein GAN (Deshpande et al., 2018). Although the maximization on discriminator’s parameters is concave, the minimization over the generator’s parameters is still non-convex due to the use of a neural-net generator. Thus, the overall optimization problem is a non-trivial non-convex concave min-max problem. From that perspective, our result in Theorem 2.1 *partially* explains the role of model overparameterization in GDA’s convergence for GANs.

Given the closed form equation 5, one may wonder what would happen if we run gradient descent on this minimization objective directly. That is running gradient descent updates of the form $\mathbf{W}_{\tau+1} = \mathbf{W}_\tau - \eta \nabla \mathcal{L}(\mathbf{W}_\tau)$ with $\mathcal{L}(\mathbf{W})$ given by equation 5. This is equivalent to GDA but instead of running one gradient ascent iteration for the maximization iteration we run infinitely many. Interestingly, in some successful GAN implementations (Gulrajani et al., 2017), often more updates on the discriminator’s parameters are run per generator’s updates. This is the subject of the next result.

Theorem 2.2 Consider the setup of Theorem 2.1. Then as long as

$$k \geq C \cdot md^4 \log(d)^3$$

²The zero initialization of \mathbf{d} is merely done for simplicity. A similar result can be derived for an arbitrary initialization of the discriminator’s parameters with minor modifications. See Theorem 2.3 for such a result.

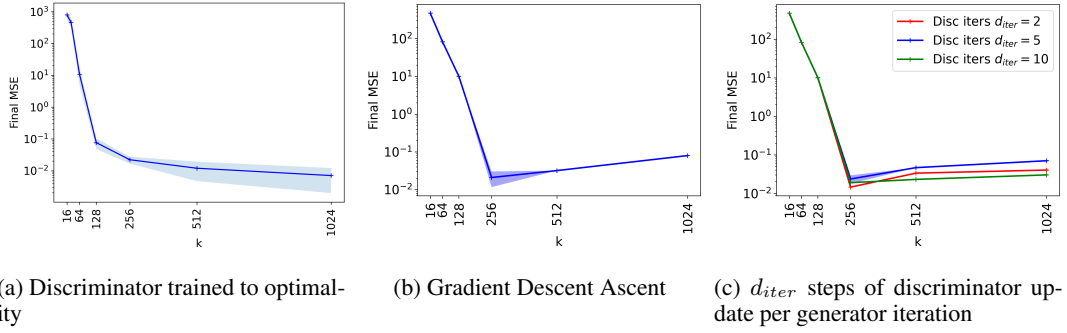


Figure 2: **Convergence plot** a GAN model with linear discriminator and 1-hidden layer generator as the hidden dimension (k) increases. *Final mse* is the mse loss between true data mean and the mean of generated distribution. Over-parameterized models show improved convergence

with C a fixed numerical constant, running GD updates of the form $\mathbf{W}_{\tau+1} = \mathbf{W}_{\tau} - \eta \nabla \mathcal{L}(\mathbf{W}_{\tau})$ on the loss given in equation 5 with step-size $\eta = \frac{2\bar{\eta}}{243k \cdot \frac{d+\frac{n-1}{\pi}}{n} \cdot \sigma_v^2 \cdot \sigma_z^2}$, with $\bar{\eta} \leq 1$, satisfies

$$\left\| \frac{1}{n} \sum_{i=1}^n \mathbf{VReLU}(\mathbf{W}_{\tau} \mathbf{z}_i) - \bar{\mathbf{x}} \right\|_{\ell_2} \leq (1 - 4 \times 10^{-6} \cdot \bar{\eta})^{\tau} \left\| \frac{1}{n} \sum_{i=1}^n \mathbf{VReLU}(\mathbf{W}_0 \mathbf{z}_i) - \bar{\mathbf{x}} \right\|_{\ell_2}. \quad (6)$$

This holds with probability at least $1 - (n+5)e^{-\frac{m}{1500}} - 5k \cdot e^{-c_1 \cdot n} - (2k+2)e^{-\frac{d}{216}} - ne^{-c_2 \cdot md^3 \log(d)^2}$ with c_1, c_2 fixed numerical constants.

This theorem states that if we solve the max part of equation 2 in closed form and run GD on the loss function per equation 5 with enough overparameterization, the loss will decrease at a geometric rate to zero. This result holds again when the model is sufficiently overparameterized. The proof of Theorem 2.2 relies on a result from (Oymak & Soltanolkotabi, 2020), which was developed in the framework of supervised learning. Also note that the amount of overparameterization required in both Theorems 2.1 and 2.2 is the same.

2.4 CAN THE ANALYSIS BE EXTENDED TO MORE GENERAL GANS?

In the previous section, we focused on the implications of our results for one-hidden layer generator and linear discriminator. However, as it will become clear in the proofs, our theoretical results are based on analyzing the convergence behavior of GDA on a more general min-max problem of the form

$$\min_{\theta \in \mathbb{R}^p} \max_{\mathbf{d} \in \mathbb{R}^m} h(\theta, \mathbf{d}) := \langle \mathbf{d}, f(\theta) - \mathbf{y} \rangle - \frac{\|\mathbf{d}\|_{\ell_2}^2}{2}, \quad (7)$$

where $f: \mathbb{R}^p \rightarrow \mathbb{R}^m$ denotes a general nonlinear mapping.

Theorem 2.3 (Informal version of Theorem A.4) Consider a general nonlinear mapping $f: \mathbb{R}^p \rightarrow \mathbb{R}^m$ with the singular values of its Jacobian mapping around initialization obeying certain assumptions (most notably $\sigma_{\min}(\mathcal{J}(\theta_0)) \geq \alpha$). Then, running GDA iterations of the form

$$\begin{cases} \mathbf{d}_{t+1} = \mathbf{d}_t + \mu \nabla_{\mathbf{d}} h(\theta_t, \mathbf{d}_t) \\ \theta_{t+1} = \theta_t - \eta \nabla_{\theta} h(\theta_t, \mathbf{d}_t) \end{cases} \quad (8)$$

with sufficiently small step sizes η and μ obeys

$$\|f(\theta_t) - \mathbf{y}\|_{\ell_2} \leq \gamma \left(1 - \frac{\eta\alpha^2}{2}\right)^t \sqrt{\|f(\theta_0) - \mathbf{y}\|_{\ell_2}^2 + \|\mathbf{d}_0\|_{\ell_2}^2}.$$

Note that similar to the previous sections one can solve the maximization problem in equation 7 in closed form so that equation 7 is equivalent to the following minimization problem

$$\min_{\theta \in \mathbb{R}^p} \mathcal{L}(\theta) := \frac{1}{2} \|f(\theta) - \mathbf{y}\|_{\ell_2}^2, \quad (9)$$

with global optima equal to zero. Theorem 2.3 ensures that GDA converges with a fast geometric rate to this global optima. This holds as soon as the model $f(\theta)$ is sufficiently overparameterized which is quantitatively captured via the minimum singular value assumption on the Jacobian at initialization ($\sigma_{\min}(\mathcal{J}(\theta_0)) \geq \alpha$ which can only hold when $m \leq p$). This general result can thus be used to provide theoretical guarantees for a much more general class of generators and discriminators. To be more specific, consider a deep GAN model where the generator \mathcal{G}_θ is a deep neural network with parameters θ and the discriminator is a deep random feature model of the form $\mathcal{D}_d(\mathbf{x}) = \mathbf{d}^T \psi(\mathbf{x})$ parameterized with d and $\psi : \mathbb{R}^d \rightarrow \mathbb{R}^m$ a deep neural network with random weights. Then the min-max training optimization problem equation 1 with a regularizer $\mathcal{R}(d) = -\|d\|_{\ell_2}^2/2$ is a special instance of equation 7 with

$$f(\theta) := \frac{1}{n} \sum_{i=1}^n \psi(\mathcal{G}_\theta(z_i)) \quad \text{and} \quad \mathbf{y} := \frac{1}{n} \sum_{i=1}^n \psi(x_i)$$

Therefore, the above result can in principle be used to rigorously analyze global convergence of GDA for an overparameterized GAN problem with a deep generator and a deep random feature discriminator model. However, characterizing the precise amount of overparameterization required for such a result to hold requires a precise analysis of the minimum singular value of the Jacobian of $f(\theta)$ at initialization as well as other singular value related conditions stated in Theorem A.4. We defer such a precise analysis to future works.

Numerical Validations: Next, we numerically study the convergence of GAN model considered in Theorems 2.1 and 2.2 where the discriminator is a linear network while the generator is a one hidden layer neural net. In our experiments, we generate x_i 's from an m -dimension Gaussian distribution with mean μ and an identity covariance matrix. The mean vector μ is randomly generated. We train two variants of GAN models using (1) GDA (as considered in Thm 2.1) and (2) GD on generator while solving the discriminator to optimality (as considered in Thm 2.2).

In Fig. 2, we plot the converged loss values of GAN models trained using both techniques (1) and (2) as the hidden dimension k of the generator is varied. The MSE loss between the true data mean and the data mean of generated samples is used as our evaluation metric. As this MSE loss approaches 0, the model converges to the global saddle point. We observe that overparameterized GAN models show improved convergence behavior than the narrower models. Additionally, the MSE loss converges to 0 for larger values of k which shows that with sufficient overparameterization, GDA converges to a global saddle point.

3 EXPERIMENTS

In this section, we demonstrate benefits of overparameterization in large GAN models. In particular, we train GANs on two benchmark datasets: CIFAR-10 (32×32 resolution) and Celeb-A (64×64 resolution). We use two commonly used GAN architectures: DCGAN and Resnet-based GAN. For both of these architectures, we train several models, each with a different number of filters in each layer, denoted by k . For simplicity, we refer to k as the hidden dimension. Appendix Fig. 8 illustrates the architectures used in our experiments. Networks with large k are more overparameterized.

We use the same value of k for both generator and discriminator networks. This is in line with the design choice made in most recent GAN models (Radford et al., 2016; Brock et al., 2019), where the size of generator and discriminator models are roughly maintained the same. We train each model till convergence and evaluate the performance of converged models using FID scores. FID scores measure the Frechet distance between feature distributions of real and generated data

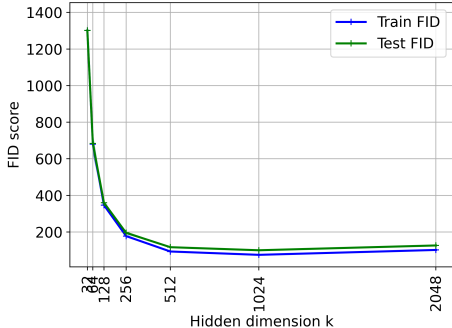


Figure 3: **MLP Overparameterization on MNIST.**

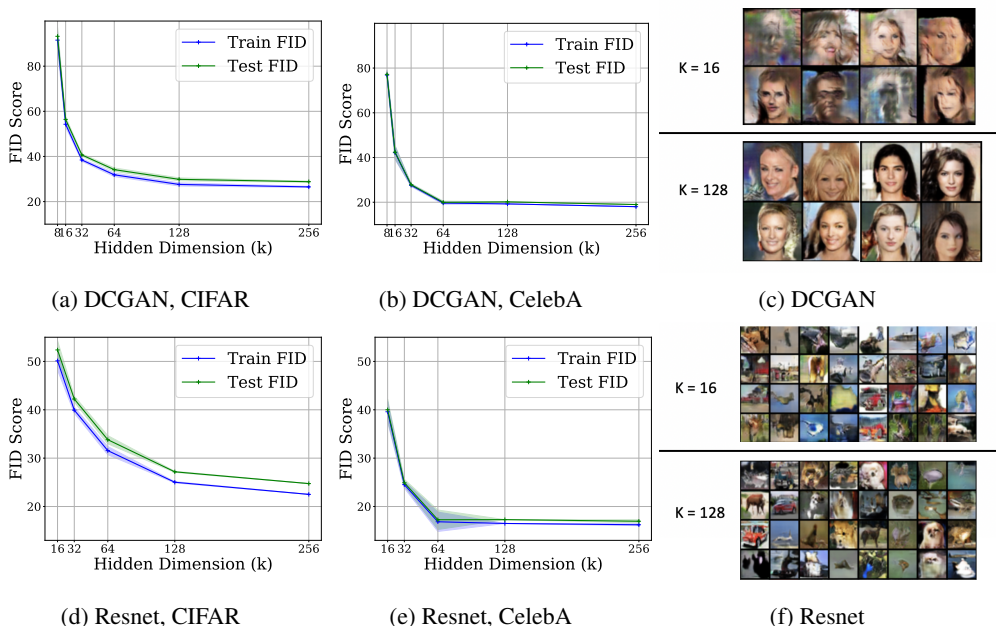


Figure 4: **Overparameterization Results:** We plot the FID scores (lower, better) of DCGAN and Resnet DCGAN as the hidden dimension k is varied. Results on CIFAR-10 and Celeb-A are shown on the plots on the left and right panels, respectively. Overparameterization gives better FID scores.

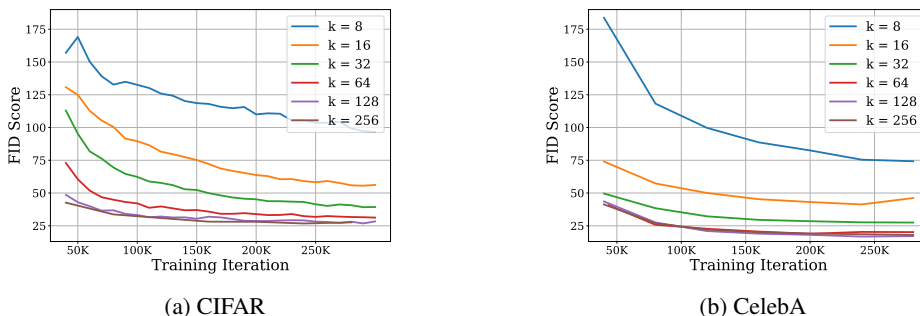


Figure 5: **DCGAN Training Results:** We plot the FID scores across training iterations of DCGAN on CIFAR-10 and Celeb-A for different values of hidden dimension k . Remarkably, we observe that over-parameterization improves the rate of convergence of GDA and its stability in training.

distributions (Heusel et al., 2017). A small FID score indicates high-quality synthesized samples. Each experiment is conducted for 5 runs, and mean and the variance of FID scores are reported.

Overparameterization yields better generative models: In Fig. 4, we show the plot of FID scores as the hidden dimension (k) is varied for DCGAN and Resnet GAN models. We observe a clear trend where the FID scores are high (i.e. poor) for small values of k , while they improve as models become more overparameterized. Also, the FID scores saturate beyond $k = 64$ for DCGAN models, and $k = 128$ for Resnet GAN models. Interestingly, these are the standard values used in the existing model architectures (Radford et al., 2016; Gulrajani et al., 2017). This trend is also consistent on MLP GANs trained on MNIST dataset (Fig. 3). We however notice that FID score in MLP GANs increase marginally as k increases from 1024 to 2048. This is potentially due to an increased generalization gap in this regime where it offsets potential benefits of over-parameterization

Overparameterization leads to improved convergence of GDA: In Fig. 5, we show the plot of FID scores over training iterations for different values for k . We observe that models with larger

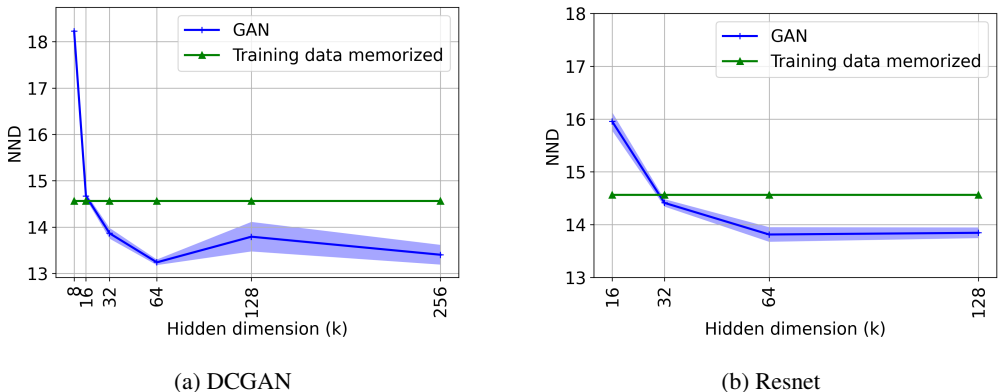


Figure 6: **Generalization in GANs:** We plot the **NND scores** as the hidden dimension k is varied for DCGAN (shown in (a)) and Resnet (shown in (b)) models.

values of k converge faster and demonstrate a more stable behavior. This agrees with our theoretical results that overparameterized models have a fast rate of convergence.

Generalization gap in GANs: To study the generalization gap, we compute the FID scores by using (1) the training-set of real data, which we call *FID train*, and (2) a held-out validation set of real data, which we call *FID test*. In Fig. 4, a plot of FID train (in blue) and FID test (in green) are shown as the hidden dimension k is varied. We observe that FID test values are consistently higher than the the FID train values. Their gap does not increase with increasing overparameterization.

However, as explained in (Gulrajani et al., 2019), the FID score has the issue of assigning low values to memorized samples. To alleviate the issue, (Gulrajani et al., 2019; Arora et al., 2017) proposed Neural Net Divergence (NND) to measure generalization in GANs. In Fig. 6, we plot NND scores by varying the hidden dimensions in DCGAN and Resnet GAN trained on CIFAR-10 dataset. We observe that increasing the value of k decreases the NND score. Interestingly, the NND score of memorized samples are higher than most of the GAN models. This indicates that overparameterized models have not been memorizing training samples and produce better generative models.

4 CONCLUSION

In this paper, we perform a systematic study of the importance of overparameterization in training GANs. We first analyze a GAN model with one-hidden layer generator and a linear discriminator optimized using Gradient Descent Ascent (GDA). Under this setup, we prove that with sufficient overparameterization, GDA converges to a global saddle point. Additionally, our result demonstrate that overparameterized models have a fast rate of convergence. We then validate our theoretical findings through extensive experiments on DCGAN and Resnet models trained on CIFAR-10 and Celeb-A datasets. We observe overparameterized models to perform well both in terms of the rate of convergence and the quality of generated samples.

5 ACKNOWLEDGEMENT

M. Sajedi would like to thank Sarah Dean for introducing (Rugh, 1996). This project was supported in part by NSF CAREER AWARD 1942230, HR00112090132, HR001119S0026, NIST 60NANB20D134 and Simons Fellowship on “Foundations of Deep Learning.” M. Soltanolkotabi is supported by the Packard Fellowship in Science and Engineering, a Sloan Research Fellowship in Mathematics, an NSF-CAREER under award 1846369, the Air Force Office of Scientific Research Young Investigator Program (AFOSR-YIP) under award FA9550-18-1-0078, DARPA Learning with Less Labels (LwLL) and FastNICS programs, and NSF-CIF awards 1813877 and 2008443.

REFERENCES

- Leonard Adolphs, Hadi Daneshmand, Aurelien Lucchi, and Thomas Hofmann. Local saddle point optimization: A curvature exploitation approach. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pp. 486–495, 2019.
- Zeyuan Allen-Zhu, Yuanzhi Li, and Zhao Song. A convergence theory for deep learning via over-parameterization. In Kamalika Chaudhuri and Ruslan Salakhutdinov (eds.), *Proceedings of the 36th International Conference on Machine Learning, ICML 2019, Long Beach, California, USA, June 9-15, 2019*, volume 97 of *Proceedings of Machine Learning Research*, pp. 242–252, 2019.
- Wangpeng An, Haoqian Wang, Qingyun Sun, Jun Xu, Qionghai Dai, and Lei Zhang. A pid controller approach for stochastic optimization of deep networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 8522–8531, 2018.
- Martin Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein generative adversarial networks. In *Proceedings of the 34th International Conference on Machine Learning, ICML 2017, Sydney, Australia, Aug 6-11, 2017*, *Proceedings of Machine Learning Research*, pp. 214–223, 2017.
- Sanjeev Arora, Rong Ge, Yingyu Liang, Tengyu Ma, and Yi Zhang. Generalization and equilibrium in generative adversarial nets (gans). *arXiv preprint arXiv:1703.00573*, 2017.
- Kenneth J Arrow, Leonid Hurwicz, and Hirofumi Uzawa. *Studies in linear and non-linear programming*. Cambridge Univ. Press, 1958.
- Hugo Berard, Gauthier Gidel, Amjad Almahairi, Pascal Vincent, and Simon Lacoste-Julien. A closer look at the optimization landscapes of generative adversarial networks. In *International Conference on Learning Representations*, 2020. URL <https://openreview.net/forum?id=HJeVnCEKwH>.
- Andrew Brock, Jeff Donahue, and Karen Simonyan. Large scale GAN training for high fidelity natural image synthesis. In *International Conference on Learning Representations*, 2019.
- Lenaic Chizat, Edouard Oyallon, and Francis Bach. On lazy training in differentiable programming. In *Advances in Neural Information Processing Systems*, pp. 2937–2947, 2019.
- Aidan Clark, Jeff Donahue, and Karen Simonyan. Adversarial video generation on complex datasets. *arXiv preprint arXiv:1907.06571*, 2019. URL <https://arxiv.org/abs/1907.06571>.
- Constantinos Daskalakis, Stratis Skoulakis, and Manolis Zampetakis. The complexity of constrained min-max optimization. *arXiv preprint arXiv:2009.09623*, 2020. URL <https://arxiv.org/abs/2009.09623>.
- Ishan Deshpande, Ziyu Zhang, and Alexander G Schwing. Generative modeling using the sliced wasserstein distance. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 3483–3491, 2018.
- Laurent Dinh, Jascha Sohl-Dickstein, and Samy Bengio. Density estimation using real nvp. In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*, 2017.
- Simon S. Du, Xiyu Zhai, Barnabas Poczos, and Aarti Singh. Gradient descent provably optimizes over-parameterized neural networks. In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019, Conference Track Proceedings*, 2019.
- Farzan Farnia and Asuman Ozdaglar. Gans may have no nash equilibria. *arXiv preprint arXiv:2002.09124*, 2020.
- Soheil Feizi, Farzan Farnia, Tony Ginart, and David Tse. Understanding GANs: the LQG setting. *arXiv preprint arXiv:1710.10793*, 2018. URL <https://arxiv.org/abs/1710.10793>.
- Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In Z. Ghahramani, M. Welling, C. Cortes, N. D. Lawrence, and K. Q. Weinberger (eds.), *Advances in Neural Information Processing Systems 27*, pp. 2672–2680, 2014.

- Ishaan Gulrajani, Faruk Ahmed, Martin Arjovsky, Vincent Dumoulin, and Aaron C Courville. Improved training of wasserstein gans. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (eds.), *Advances in Neural Information Processing Systems 30*, pp. 5767–5777, 2017.
- Ishaan Gulrajani, Colin Raffel, and Luke Metz. Towards GAN benchmarks which require generalization. In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=HkxKH2AcFm>.
- Shaul Gutman. Uncertain dynamical systems—a lyapunov min-max approach. *IEEE Transactions on Automatic Control*, 24(3):437–443, 1979.
- Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. In *Advances in neural information processing systems 30*, pp. 6626–6637, 2017.
- Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2019, Long Beach, CA, USA, June 16-20, 2019*, pp. 4401–4410. Computer Vision Foundation / IEEE, 2019.
- Diederik P. Kingma and Max Welling. Auto-encoding variational bayes. In Yoshua Bengio and Yann LeCun (eds.), *2nd International Conference on Learning Representations, ICLR 2014, Banff, AB, Canada, April 14-16, 2014, Conference Track Proceedings*, 2014.
- Michel Ledoux. *The concentration of measure phenomenon*. Number 89 in Mathematical Surveys and Monographs. American Mathematical Soc., 2001.
- Qi Lei, Jason D Lee, Alexandros G Dimakis, and Constantinos Daskalakis. Sgd learns one-layer networks in wgens. *arXiv preprint arXiv:1910.07030*, 2019. URL <https://arxiv.org/abs/1910.07030>.
- Chun-Liang Li, Wei-Cheng Chang, Yu Cheng, Yiming Yang, and Barnabas Poczos. Mmd gan: Towards deeper understanding of moment matching network. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (eds.), *Advances in Neural Information Processing Systems 30*, pp. 2203–2213, 2017.
- Lars Mescheder, Sebastian Nowozin, and Andreas Geiger. The numerics of GANs. In *Advances in Neural Information Processing Systems 30*, pp. 1823–1833, 2017.
- Lars Mescheder, Andreas Geiger, and Sebastian Nowozin. Which training methods for gans do actually converge? *arXiv preprint arXiv:1801.04406*, 2018. URL <https://arxiv.org/abs/1801.04406>.
- Vaishnavh Nagarajan and J Zico Kolter. Gradient descent GAN optimization is locally stable. In *Advances in Neural Information Processing Systems 30*, pp. 5591–5600, 2017.
- Samet Oymak and Mahdi Soltanolkotabi. Overparameterized nonlinear learning: Gradient descent takes the shortest path? In *Proceedings of the 36th International Conference on Machine Learning, ICML 2019, Long Beach, California, USA, June 9-15, 2019*, Proceedings of Machine Learning Research, pp. 4951–4960, 2019.
- Samet Oymak and Mahdi Soltanolkotabi. Towards moderate overparameterization: global convergence guarantees for training shallow neural networks. *IEEE Journal on Selected Areas in Information Theory*, 2020.
- Samet Oymak, Zalan Fabian, Mingchen Li, and Mahdi Soltanolkotabi. Generalization guarantees for neural networks via harnessing the low-rank structure of the jacobian. *arXiv preprint arXiv:1906.05392*, 2019. URL <https://arxiv.org/abs/1906.05392>.
- Alec Radford, Luke Metz, and Soumith Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. In Yoshua Bengio and Yann LeCun (eds.), *4th International Conference on Learning Representations, ICLR 2016, San Juan, Puerto Rico, May 2-4, 2016, Conference Track Proceedings*, 2016.

- Wilson J Rugh. *Linear system theory*. Prentice-Hall, Inc., 1996.
- Mahdi Soltanolkotabi. Structured signal recovery from quadratic measurements: Breaking sample complexity barriers via nonconvex optimization. *IEEE Transactions on Information Theory*, 65(4):2374–2400, 2019.
- Mahdi Soltanolkotabi, Adel Javanmard, and Jason D Lee. Theoretical insights into the optimization landscape of over-parameterized shallow neural networks. *IEEE Transactions on Information Theory*, 65(2):742–769, 2018.
- Dave Van Veen, Ajil Jalal, Mahdi Soltanolkotabi, Eric Price, Sriram Vishwanath, and Alexandros G Dimakis. Compressed sensing with deep image prior and learned regularization. *arXiv preprint arXiv:1806.06438*, 2018. URL <https://arxiv.org/abs/1806.06438>.
- John Von Neumann and Oskar Morgenstern. *Theory of games and economic behavior*. Princeton university press, 1953.
- Kun Xu, Chongxuan Li, Huanshu Wei, Jun Zhu, and Bo Zhang. Understanding and stabilizing gans’ training dynamics with control theory. *arXiv preprint arXiv:1909.13188*, 2019.
- Difan Zou and Quanquan Gu. An improved analysis of training over-parameterized deep neural networks. In *Advances in Neural Information Processing Systems 32*, pp. 2055–2064, 2019.

Appendix

A PROOFS

In this section, we prove Theorems 2.1 and 2.2. First, we provide some notations we use throughout the remainder of the paper in Section A.1. Before proving these specialized results for one hidden layer generators and linear discriminators (Theorems 2.1 and 2.2), we state and prove a more general result (formal version of Theorem 2.3) on the convergence of GDA on a general class of min-max problems in Section A.3. Then we state a few preliminary calculations in Section A.4. Next, we state some key lemmas in Section A.5 and defer their proofs to Appendix B. Finally, we prove Theorems 2.1 and 2.2 in Sections A.6 and A.7, respectively.

A.1 NOTATION

We will use C, c, c_1 , etc. to denote positive absolute constants, whose value may change throughout the paper and from line to line. We use $\phi(z) = \text{ReLU}(z) = \max(0, z)$ and its (generalized) derivative $\phi'(z) = \mathbf{1}_{\{z \geq 0\}}$ with $\mathbf{1}$ being the indicator function. $\sigma_{\min}(\mathbf{X})$ and $\sigma_{\max}(\mathbf{X}) = \|\mathbf{X}\|$ denote the minimum and maximum singular values of matrix \mathbf{X} . For two arbitrary matrices \mathbf{A} and \mathbf{B} , $\mathbf{A} \otimes \mathbf{B}$ denotes their kronecker product. The spectral radius of a matrix $\mathbf{A} \in \mathbb{C}^{n \times n}$ is defined as $\rho(\mathbf{A}) = \max\{|\lambda_1|, \dots, |\lambda_n|\}$, where λ_i 's are the eigenvalues of \mathbf{A} . Throughout the proof we shall assume $\phi := \text{ReLU}$ to avoid unnecessarily long expressions.

A.2 PROOF SKETCH OF THE MAIN RESULTS

In this section, we provide a brief overview of our proofs. We focus on the main result in this manuscript, which is about the convergence of GDA (Theorem 2.1). To do this we study the convergence of GDA on the more general min-max problem of the form (see Theorem A.4 for a formal statement)

$$\min_{\boldsymbol{\theta} \in \mathbb{R}^n} \max_{\mathbf{d} \in \mathbb{R}^m} h(\boldsymbol{\theta}, \mathbf{d}) := \langle \mathbf{d}, f(\boldsymbol{\theta}) - \mathbf{y} \rangle - \frac{\|\mathbf{d}\|_{\ell_2}^2}{2}. \quad (10)$$

In this case the GDA iterates take the form

$$\begin{cases} \mathbf{d}_{t+1} = (1 - \mu) \mathbf{d}_t + \mu (f(\boldsymbol{\theta}_t) - \mathbf{y}) \\ \boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \eta \mathcal{J}^T(\boldsymbol{\theta}_t) \mathbf{d}_t \end{cases}. \quad (11)$$

Our proof for global convergence of GDA on this min-max loss consists of the following steps.

Step 1: Recasting the GDA updates as a linear time-varying system

In the first step we carry out a series of algebraic manipulations to recast the GDA updates (equation 11) into the following form

$$\begin{bmatrix} \mathbf{r}_{t+1} \\ \mathbf{d}_{t+1} \end{bmatrix} = \mathbf{A}_t \begin{bmatrix} \mathbf{r}_t \\ \mathbf{d}_t \end{bmatrix},$$

where $\mathbf{r}_t = f(\boldsymbol{\theta}_t) - \mathbf{y}$ denotes the residuum and \mathbf{A}_t denotes a properly defined transition matrix.

Step 2: Approximation by a linear time-invariant system

Next, to analyze the behavior of the time-varying dynamical system above we approximate it by the following time-invariant linear dynamical system

$$\begin{bmatrix} \tilde{\mathbf{r}}_{t+1} \\ \tilde{\mathbf{d}}_{t+1} \end{bmatrix} = \begin{bmatrix} \mathbf{I} & -\eta \mathcal{J}^T(\boldsymbol{\theta}_0) \mathcal{J}(\boldsymbol{\theta}_0) \\ \mu \mathbf{I} & (1 - \mu) \mathbf{I} \end{bmatrix} \begin{bmatrix} \tilde{\mathbf{r}}_t \\ \tilde{\mathbf{d}}_t \end{bmatrix},$$

where $\boldsymbol{\theta}_0$ denotes the initialization. The validity of this approximation is ensured by our assumptions on the Jacobian of the function f , which, among others, guarantee that it does not change too much in a sufficiently large neighborhood around the initialization and that the smallest singular value of $\mathcal{J}(\boldsymbol{\theta}_0)$ is bounded from below.

Step 3: Analysis of time-invariant linear dynamical system

To analyze the time-invariant dynamical system above, we utilize and refine intricate arguments

from the control theory literature involving the spectral radius of the fixed transition matrix above to obtain

$$\left\| \begin{bmatrix} \tilde{\mathbf{r}}_t \\ \tilde{\mathbf{d}}_t \end{bmatrix} \right\|_{\ell_2} \lesssim (1 - \eta\alpha^2)^t \left\| \begin{bmatrix} \tilde{\mathbf{r}}_0 \\ \tilde{\mathbf{d}}_0 \end{bmatrix} \right\|_{\ell_2}.$$

Step 4: Completing the proof via a perturbation argument

In the last step of our proof we show that the two sequences $\begin{bmatrix} \mathbf{r}_t \\ \mathbf{d}_t \end{bmatrix}$ and $\begin{bmatrix} \tilde{\mathbf{r}}_t \\ \tilde{\mathbf{d}}_t \end{bmatrix}$ will remain close to each other. This is based on a novel perturbation argument. The latter combined with Step 3 allows us to conclude

$$\left\| \begin{bmatrix} \mathbf{r}_t \\ \mathbf{d}_t \end{bmatrix} \right\|_{\ell_2} \lesssim \left(1 - \frac{\eta\alpha^2}{2}\right)^t \left\| \begin{bmatrix} \mathbf{r}_0 \\ \mathbf{d}_0 \end{bmatrix} \right\|_{\ell_2},$$

which finishes the global convergence of GDA on equation 10 and hence the proof of Theorem A.4.

In order to deduce Theorem 2.1 from Theorem A.4, we need to check that the Jacobian at the initialization is bounded from below at the origin and that it does not change too quickly in a large enough neighborhood. In order to prove that we will leverage recent ideas from the deep learning theory literature revolving around the neural tangent kernel. This allows us to guarantee that this conditions are indeed met, if the neural network is sufficiently wide and the initialization is chosen large enough.

The second main result of this manuscript, Theorem 2.2, can be deduced more directly from recent results on overparameterized learning (see Oymak & Soltanolkotabi (2020)). Hence, we have deferred its proof to Section A.7.

A.3 ANALYSIS OF GDA: A CONTROL THEORY PERSPECTIVE

In this section we will focus on solving a general min-max optimization problem of the form

$$\min_{\boldsymbol{\theta} \in \mathbb{R}^n} \max_{\mathbf{d} \in \mathbb{R}^m} h(\boldsymbol{\theta}, \mathbf{d}) := \langle \mathbf{d}, f(\boldsymbol{\theta}) - \mathbf{y} \rangle - \frac{\|\mathbf{d}\|_{\ell_2}^2}{2}, \quad (12)$$

where $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is a general nonlinear mapping. In particular, we focus on analyzing the convergence behavior of Gradient Descent/Ascent (GDA) on the above loss, starting from initial estimates $\boldsymbol{\theta}_0$ and \mathbf{d}_0 . In this case the GDA updates take the following form

$$\begin{cases} \mathbf{d}_{t+1} = (1 - \mu) \mathbf{d}_t + \mu (f(\boldsymbol{\theta}_t) - \mathbf{y}) \\ \boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \eta \mathcal{J}^T(\boldsymbol{\theta}_t) \mathbf{d}_t \end{cases}. \quad (13)$$

We note that solving the inner maximization problem in equation 12 would yield

$$\min_{\boldsymbol{\theta} \in \mathbb{R}^n} \frac{1}{2} \|f(\boldsymbol{\theta}) - \mathbf{y}\|_{\ell_2}^2. \quad (14)$$

In this section, our goal is to show that when running the GDA updates of equation 13, the norm of the residual vector defined as $\mathbf{r}_t := f(\boldsymbol{\theta}_t) - \mathbf{y}$ goes to zero and hence we reach a global optimum of equation 14 (and in turn equation 12).

Our proof will build on ideas from control theory and dynamical systems literature. For that, we are first going to rewrite the equations 13 in a more convenient way. We define the average Jacobian along the path connecting two points $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ as

$$\mathcal{J}(\mathbf{y}, \mathbf{x}) = \int_0^1 \mathcal{J}(\mathbf{x} + \alpha(\mathbf{y} - \mathbf{x})) d\alpha,$$

where $\mathcal{J}(\boldsymbol{\theta}) \in \mathbb{R}^{m \times n}$ is the Jacobian associated with the nonlinear mapping f . Next, from the fundamental theorem of calculus it follows that

$$\begin{aligned} \mathbf{r}_{t+1} &= f(\boldsymbol{\theta}_{t+1}) - \mathbf{y} = f(\boldsymbol{\theta}_t - \eta \mathcal{J}_t^T \mathbf{d}_t) - \mathbf{y} \\ &= f(\boldsymbol{\theta}_t) - \eta \mathcal{J}_{t+1,t} \mathcal{J}_t^T \mathbf{d}_t - \mathbf{y} \\ &= \mathbf{r}_t - \eta \mathcal{J}_{t+1,t} \mathcal{J}_t^T \mathbf{d}_t, \end{aligned} \quad (15)$$

where we used the shorthands $\mathcal{J}_t := \mathcal{J}(\boldsymbol{\theta}_t)$ and $\mathcal{J}_{t+1,t} := \mathcal{J}(\boldsymbol{\theta}_{t+1}, \boldsymbol{\theta}_t)$ for exposition purposes.

Next, we combine the updates \mathbf{r}_t and \mathbf{d}_t into a state vector of the form $\mathbf{z}_t := \begin{bmatrix} \mathbf{r}_t \\ \mathbf{d}_t \end{bmatrix} \in \mathbb{R}^{2m}$. Using this notation the relationship between the state vectors from one iteration to the next takes the form

$$\mathbf{z}_{t+1} = \underbrace{\begin{bmatrix} \mathbf{I} & -\eta \mathcal{J}_{t+1,t} \mathcal{J}_t^T \\ \mu \mathbf{I} & (1-\mu) \mathbf{I} \end{bmatrix}}_{=: \mathbf{A}_t} \mathbf{z}_t, \quad t \geq 0, \quad (16)$$

which resembles a time-varying linear dynamical system with transition matrix \mathbf{A}_t . Now note that to show convergence of \mathbf{r}_t to zero it suffices to show convergence of \mathbf{z}_t to zero. To do this we utilize the following notion of uniform exponential stability, which will be crucial in analyzing the solutions of equation 16. (See [Rugh \(1996\)](#) for a comprehensive overview on stability notions in discrete state equations.)

Definition 1 *A linear state equation of the form $\mathbf{z}_{t+1} = \mathbf{A}_t \mathbf{z}_t$ is called uniformly exponentially stable if for every $t \geq 0$ we have $\|\mathbf{z}_t\|_{\ell_2} \leq \gamma \lambda^t \|\mathbf{z}_0\|_{\ell_2}$, where $\gamma \geq 1$ is a finite constant and $0 \leq \lambda < 1$.*

Using the above definition to show the convergence of the state vector \mathbf{z}_t to zero at a geometric rate it suffices to show the state equation 16 is exponentially stable.³ For that, we are first going to analyze a state equation which results from linearizing the nonlinear function $f(\boldsymbol{\theta})$ around the initialization $\boldsymbol{\theta}_0$. In the next step, we are going to show that the behavior of these two problems are similar, provided we stay close to initialization (which we are also going to prove). Specifically, we consider the linearized problem

$$\min_{\tilde{\boldsymbol{\theta}} \in \mathbb{R}^n} \max_{\tilde{\mathbf{d}} \in \mathbb{R}^m} h_{\text{lin}}(\tilde{\boldsymbol{\theta}}, \tilde{\mathbf{d}}) := \left\langle \tilde{\mathbf{d}}, f(\boldsymbol{\theta}_0) + \mathcal{J}_0(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}_0) - \mathbf{y} \right\rangle - \frac{\|\tilde{\mathbf{d}}\|_{\ell_2}^2}{2}. \quad (17)$$

We first analyze GDA on this linearized problem starting from the same initialization as the original problem, i.e. $\tilde{\boldsymbol{\theta}}_0 = \boldsymbol{\theta}_0$ and $\tilde{\mathbf{d}}_0 = \mathbf{d}_0$. The gradient descent update for $\tilde{\boldsymbol{\theta}}_t$ takes the form

$$\tilde{\boldsymbol{\theta}}_{t+1} = \tilde{\boldsymbol{\theta}}_t - \eta \mathcal{J}_0^T \tilde{\mathbf{d}}_t, \quad (18)$$

and the gradient ascent update for $\tilde{\mathbf{d}}_t$ takes the form

$$\begin{aligned} \tilde{\mathbf{d}}_{t+1} &= \tilde{\mathbf{d}}_t + \mu \left(f(\boldsymbol{\theta}_0) + \mathcal{J}_0(\tilde{\boldsymbol{\theta}}_t - \boldsymbol{\theta}_0) - \mathbf{y} - \tilde{\mathbf{d}}_t \right) \\ &= (1-\mu) \tilde{\mathbf{d}}_t + \mu \tilde{\mathbf{r}}_t, \end{aligned} \quad (19)$$

where we used the linear residual defined as $\tilde{\mathbf{r}}_t = f(\boldsymbol{\theta}_0) + \mathcal{J}_0(\tilde{\boldsymbol{\theta}}_t - \boldsymbol{\theta}_0) - \mathbf{y}$. Moreover, the residual from one iterate to the next can be written as follows

$$\begin{aligned} \tilde{\mathbf{r}}_{t+1} &= f(\boldsymbol{\theta}_0) + \mathcal{J}_0(\tilde{\boldsymbol{\theta}}_{t+1} - \boldsymbol{\theta}_0) - \mathbf{y} \\ &= f(\boldsymbol{\theta}_0) + \mathcal{J}_0(\tilde{\boldsymbol{\theta}}_t - \eta \mathcal{J}_0^T \tilde{\mathbf{d}}_t - \boldsymbol{\theta}_0) - \mathbf{y} \\ &= \tilde{\mathbf{r}}_t - \eta \mathcal{J}_0 \mathcal{J}_0^T \tilde{\mathbf{d}}_t. \end{aligned} \quad (20)$$

Again, we define a new vector $\tilde{\mathbf{z}}_t = \begin{bmatrix} \tilde{\mathbf{r}}_t \\ \tilde{\mathbf{d}}_t \end{bmatrix} \in \mathbb{R}^{2m}$ and by putting together equations 19 and 20 we arrive at

$$\tilde{\mathbf{z}}_{t+1} = \begin{bmatrix} \mathbf{I} & -\eta \mathcal{J}_0 \mathcal{J}_0^T \\ \mu \mathbf{I} & (1-\mu) \mathbf{I} \end{bmatrix} \tilde{\mathbf{z}}_t = \mathbf{A} \tilde{\mathbf{z}}_t, \quad t \geq 0, \quad (21)$$

³We note that technically the dynamical system equation 16 is not linear. However, we still use exponential stability with some abuse of notation to refer to the property that $\|\mathbf{z}_t\|_{\ell_2} \leq \gamma \lambda^t \|\mathbf{z}_0\|_{\ell_2}$ holds. As we will see in the forth-coming paragraphs, our formal analysis is via a novel perturbation analysis of a linear dynamical system and therefore keeping this terminology is justified.

which is of the form of a linear time-invariant state equation. As a first step in our proof, we are going to show that the linearized state equations are uniformly exponentially stable. First, recall the following well-known lemma, which characterizes uniform exponential stability in terms of the eigenvalues of \mathbf{A} .

Lemma A.1 (*Rugh, 1996, Theorem 22.11*) *A linear state equation of the form $\tilde{\mathbf{z}}_{t+1} = \mathbf{A}\tilde{\mathbf{z}}_t$ with \mathbf{A} a fixed matrix is uniformly exponentially stable if and only if all eigenvalues of \mathbf{A} have magnitudes strictly less than one, i.e. $\rho(\mathbf{A}) < 1$. In this case, it holds for all $t \geq 0$ and all \mathbf{z} that*

$$\|\mathbf{A}^t \mathbf{z}\| \leq \gamma \rho(\mathbf{A})^t \|\mathbf{z}\|,$$

where $\gamma \geq 1$ is an absolute constant, which only depends on \mathbf{A} .

In the next lemma, we prove that under suitable assumptions on \mathcal{J}_0 and the step sizes μ and η the state equations 21 indeed fulfill this condition.

Lemma A.2 *Assume that $\alpha \leq \sigma_{\min}(\mathcal{J}_0) \leq \sigma_{\max}(\mathcal{J}_0) \leq \beta$ and consider the matrix $\mathbf{A} = \begin{bmatrix} \mathbf{I} & -\eta \mathcal{J}_0 \mathcal{J}_0^T \\ \mu \mathbf{I} & (1-\mu) \mathbf{I} \end{bmatrix}$. Suppose that $\frac{\mu}{\eta} \geq 4\beta^2$. Then it holds that $\rho(\mathbf{A}) \leq 1 - \eta\alpha^2$.*

Proof Suppose that λ is an eigenvalue of \mathbf{A} . Hence, there is an eigenvector $[\mathbf{x}, \mathbf{y}]^T \neq \mathbf{0}$ such that

$$\begin{bmatrix} \mathbf{I} & -\eta \mathcal{J}_0 \mathcal{J}_0^T \\ \mu \mathbf{I} & (1-\mu) \mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} = \lambda \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix}$$

holds. By a direct calculation we observe that this yields the equation

$$\eta \mathcal{J}_0 \mathcal{J}_0^T \mathbf{x} = \left(\frac{-(1-\lambda)^2}{\mu} + (1-\lambda) \right) \mathbf{x}.$$

In particular, \mathbf{x} must be an eigenvector of $\mathcal{J}_0 \mathcal{J}_0^T$. Denoting the corresponding eigenvalue with s , we obtain the identity

$$\frac{(1-\lambda)^2}{\mu} - (1-\lambda) + \eta s = 0.$$

Hence, we must have

$$\lambda \in \left\{ 1 - \frac{\mu}{2} + \sqrt{\frac{\mu^2}{4} - \mu\eta s}; 1 - \frac{\mu}{2} - \sqrt{\frac{\mu^2}{4} - \mu\eta s} \right\}.$$

Note that the square root is indeed well-defined, since

$$\frac{\mu^2}{4} - \mu\eta s \geq \mu\eta\beta^2 - \mu\eta s \geq 0,$$

where in the first inequality we used the assumption $\frac{\mu}{\eta} \geq 4\beta^2$ and in the second line we used that $s \leq \beta^2$, which is a consequence of our assumption on the singular values of \mathcal{J}_0 . Hence, it follows by the reverse triangle inequality that

$$|\lambda| - \left(1 - \frac{\mu}{2}\right) \leq \left| \lambda - \left(1 - \frac{\mu}{2}\right) \right| = \sqrt{\left(\frac{\mu}{2}\right)^2 - \mu\eta s} < \frac{\mu}{2} - \eta s \leq \frac{\mu}{2} - \eta\alpha^2,$$

where the second inequality is valid as $\frac{\mu}{2} - \eta s \geq 0$ is implied by $\frac{\mu}{2} \geq 2\eta\beta^2 > \eta s$. In the last inequality we used the fact that $\alpha^2 \leq s$, which is a consequence of our assumption on the singular values of \mathcal{J}_0 . By rearranging terms, we obtain that $|\lambda| < 1 - \eta\alpha^2$. Since λ was an arbitrary eigenvalue of \mathbf{A} , the result follows. ■

Since the last lemma shows that under suitable conditions it holds that $\rho(\mathbf{A}) < 1$, Lemma A.3 yields uniform exponential stability of our state equations. However, this will not be sufficient for our purposes. The reason is that Lemma A.3 does not specify the constant γ and that in order to deal with the time-varying dynamical system we will need a precise estimate. The next lemma shows that for the state equations 21 we have, under suitable assumptions, $\gamma \leq 5$.

Lemma A.3 Consider the linear, time invariant system of equations

$$\tilde{\mathbf{z}}_{t+1} = \begin{bmatrix} \mathbf{I} & -\eta \mathcal{J}_0 \mathcal{J}_0^T \\ \mu \mathbf{I} & (1-\mu) \mathbf{I} \end{bmatrix} \tilde{\mathbf{z}}_t = \mathbf{A} \tilde{\mathbf{z}}_t, \quad t \geq 0.$$

Furthermore, assume that $\alpha \leq \sigma_{\min}(\mathcal{J}_0) \leq \sigma_{\max}(\mathcal{J}_0) \leq \beta$ and suppose that the condition $\frac{\mu}{\eta} \geq 8\beta^2$ is satisfied. Then there is a constant $\gamma \leq 5$ such that for all $t \geq 0$ it holds that

$$\|\tilde{\mathbf{z}}_t\|_{\ell_2} \leq \gamma (1 - \eta\alpha^2)^t \|\tilde{\mathbf{z}}_0\|_{\ell_2}.$$

Proof Denote the SVD decomposition of \mathcal{J}_0 by $\mathbf{W}\Sigma\mathbf{V}^T$ and note that

$$\begin{bmatrix} \mathbf{I} & -\eta \mathcal{J}_0 \mathcal{J}_0^T \\ \mu \mathbf{I} & (1-\mu) \mathbf{I} \end{bmatrix} = \begin{bmatrix} \mathbf{W} & \mathbf{0} \\ \mathbf{0} & \mathbf{W} \end{bmatrix} \begin{bmatrix} \mathbf{I} & -\eta \Sigma \Sigma^T \\ \mu \mathbf{I} & (1-\mu) \mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{W}^T & \mathbf{0} \\ \mathbf{0} & \mathbf{W}^T \end{bmatrix}.$$

This means we can write

$$\begin{bmatrix} \mathbf{I} & -\eta \mathcal{J}_0 \mathcal{J}_0^T \\ \mu \mathbf{I} & (1-\mu) \mathbf{I} \end{bmatrix} = \begin{bmatrix} \mathbf{W} & \mathbf{0} \\ \mathbf{0} & \mathbf{W} \end{bmatrix} \mathbf{P} \begin{bmatrix} \mathbf{C}_1 & & \mathbf{0} \\ & \ddots & \\ \mathbf{0} & & \mathbf{C}_m \end{bmatrix} \mathbf{P}^T \begin{bmatrix} \mathbf{W}^T & \mathbf{0} \\ \mathbf{0} & \mathbf{W}^T \end{bmatrix},$$

where \mathbf{P} is a permutation matrix and the matrices \mathbf{C}_i are of the form $\mathbf{C}_i = \begin{bmatrix} 1 & -\eta\sigma_i^2 \\ \mu & (1-\mu) \end{bmatrix}$, for $1 \leq i \leq m$, where the σ_i 's denote the singular values of \mathcal{J}_0 . Using this decomposition we can deduce

$$\|\tilde{\mathbf{z}}_t\|_{\ell_2} = \|\mathbf{A}^t \tilde{\mathbf{z}}_0\|_{\ell_2} \leq \|\mathbf{A}^t\| \|\tilde{\mathbf{z}}_0\|_{\ell_2} = \left(\max_{1 \leq i \leq m} \|\mathbf{C}_i^t\| \right) \|\tilde{\mathbf{z}}_0\|_{\ell_2}.$$

Now suppose that $\mathbf{V}_i \mathbf{D}_i \mathbf{V}_i^{-1}$ is the eigenvalue decomposition of \mathbf{C}_i , where the columns of \mathbf{V}_i contain the eigenvectors and \mathbf{D}_i is a diagonal matrix consisting of the eigenvalues. (Note that it follows from our assumptions on μ and η that the matrix \mathbf{C}_i is diagonalizable.) We have

$$\|\mathbf{C}_i^t\| = \|\mathbf{V}_i \mathbf{D}_i^t \mathbf{V}_i^{-1}\| \leq \|\mathbf{V}_i\| \|\mathbf{D}_i^t\| \|\mathbf{V}_i^{-1}\| = \kappa_i \cdot \rho(\mathbf{C}_i)^t,$$

where we defined $\kappa_i := \|\mathbf{V}_i\| \|\mathbf{V}_i^{-1}\|$. From Lemma A.2 we know that the assumption $\frac{\mu}{\eta} \geq 4\beta^2$ results in $\rho(\mathbf{A}) \leq 1 - \eta\alpha^2$. Therefore, defining $\gamma := \max_{1 \leq i \leq m} \kappa_i$ and noting $\rho(\mathbf{A}) = \max_{1 \leq i \leq m} \rho(\mathbf{C}_i)$, we obtain that

$$\|\tilde{\mathbf{z}}_t\|_{\ell_2} \leq \left(\max_{1 \leq i \leq m} \|\mathbf{C}_i^t\| \right) \|\tilde{\mathbf{z}}_0\|_{\ell_2} \leq \gamma (1 - \eta\alpha^2)^t \|\tilde{\mathbf{z}}_0\|_{\ell_2}.$$

In order to finish the proof we need to show that $\gamma \leq 5$. For that, note that calculating the eigenvectors of \mathbf{C}_i directly reveals that we can represent this matrix as

$$\mathbf{V}_i = \begin{bmatrix} \frac{1 + \sqrt{1 - 4\frac{\eta\sigma_i^2}{\mu}}}{2} & \frac{1 - \sqrt{1 - 4\frac{\eta\sigma_i^2}{\mu}}}{2} \\ \frac{1}{1} & \frac{1}{1} \end{bmatrix}.$$

Since $\|\mathbf{V}_i\| = \sqrt{\lambda_{\max}(\mathbf{V}_i \mathbf{V}_i^T)}$ and $\|\mathbf{V}_i^{-1}\| = \sqrt{\lambda_{\min}(\mathbf{V}_i \mathbf{V}_i^T)}$, we calculate $\mathbf{V}_i \mathbf{V}_i^T$, which yields

$$\mathbf{V}_i \mathbf{V}_i^T = \begin{bmatrix} 1 - 2\frac{\eta\sigma_i^2}{\mu} & 1 \\ 1 & 2 \end{bmatrix}.$$

This representation allows us to directly calculate the two eigenvalues of $\mathbf{V}_i \mathbf{V}_i^T$, which shows that

$$\begin{aligned} \kappa_i &= \sqrt{\frac{\lambda_{\max}(\mathbf{V}_i \mathbf{V}_i^T)}{\lambda_{\min}(\mathbf{V}_i \mathbf{V}_i^T)}} \\ &= \sqrt{\frac{3 - 2\frac{\eta\sigma_i^2}{\mu} + \sqrt{\left(1 + 2\frac{\eta\sigma_i^2}{\mu}\right)^2 + 4}}{3 - 2\frac{\eta\sigma_i^2}{\mu} - \sqrt{\left(1 + 2\frac{\eta\sigma_i^2}{\mu}\right)^2 + 4}}} \\ &= \frac{3 - 2\frac{\eta\sigma_i^2}{\mu} + \sqrt{\left(1 + 2\frac{\eta\sigma_i^2}{\mu}\right)^2 + 4}}{2\sqrt{1 - 4\frac{\eta\sigma_i^2}{\mu}}} \\ &\leq \frac{6}{2\sqrt{1 - 4\frac{\eta\sigma_i^2}{\mu}}} \\ &< 5, \end{aligned}$$

where the last inequality holds because of $\frac{\eta\sigma_i^2}{\mu} \leq \frac{\eta\beta^2}{\mu} \leq \frac{1}{8}$. Since $\gamma = \max_{1 \leq i \leq m} \kappa_i$, this finishes the proof. ■

Now that we have shown that the linearized iterates converge to the global optima we turn our attention to showing that the nonlinear iterates 16 are close to its linear counterpart 21. For that, we make the following assumptions.

Assumption 1: The singular values of the Jacobian at initialization are bounded from below

$$\sigma_{\min}(\mathcal{J}(\boldsymbol{\theta}_0)) \geq \alpha$$

for some positive constants α and β .

Assumption 2: In a neighborhood with radius R around the initialization, the Jacobian mapping associated with f obeys

$$\|\mathcal{J}(\boldsymbol{\theta})\| \leq \beta$$

for all $\boldsymbol{\theta} \in \mathcal{B}_R(\boldsymbol{\theta}_0)$, where $\mathcal{B}_R(\boldsymbol{\theta}_0) := \{\boldsymbol{\theta} \in \mathbb{R}^p : \|\boldsymbol{\theta} - \boldsymbol{\theta}_0\|_{\ell_2} \leq R\}$.

Assumption 3: In a neighborhood with radius R around the initialization, the spectral norm of the Jacobian varies no more than ϵ in the sense that

$$\|\mathcal{J}(\boldsymbol{\theta}) - \mathcal{J}(\boldsymbol{\theta}_0)\| \leq \epsilon$$

for all $\boldsymbol{\theta} \in \mathcal{B}_R(\boldsymbol{\theta}_0)$.

With these assumptions in place, we are ready to state the main theorem.

Theorem A.4 Consider the GDA updates for the min-max optimization problem 12

$$\begin{bmatrix} \mathbf{d}_{t+1} \\ \boldsymbol{\theta}_{t+1} \end{bmatrix} = \begin{bmatrix} \mathbf{d}_t + \mu \nabla_{\mathbf{d}} h(\boldsymbol{\theta}_t, \mathbf{d}_t) \\ \boldsymbol{\theta}_t - \eta \nabla_{\boldsymbol{\theta}} h(\boldsymbol{\theta}_t, \mathbf{d}_t) \end{bmatrix} \quad (22)$$

and consider the GDA updates of the linearized problem 21

$$\begin{bmatrix} \tilde{\mathbf{d}}_{t+1} \\ \tilde{\boldsymbol{\theta}}_{t+1} \end{bmatrix} = \begin{bmatrix} \tilde{\mathbf{d}}_t + \mu \nabla_{\mathbf{d}} h_{lin}(\tilde{\boldsymbol{\theta}}_t, \tilde{\mathbf{d}}_t) \\ \tilde{\boldsymbol{\theta}}_t - \eta \nabla_{\boldsymbol{\theta}} h_{lin}(\tilde{\boldsymbol{\theta}}_t, \tilde{\mathbf{d}}_t) \end{bmatrix}. \quad (23)$$

Set $\mathbf{z}_t := \begin{bmatrix} \mathbf{r}_t \\ \mathbf{d}_t \end{bmatrix}$ and $\tilde{\mathbf{z}}_t := \begin{bmatrix} \tilde{\mathbf{r}}_t \\ \tilde{\mathbf{d}}_t \end{bmatrix}$, where $\mathbf{r}_t := f(\boldsymbol{\theta}_t) - \mathbf{y}$ and $\tilde{\mathbf{r}}_t = f(\boldsymbol{\theta}_0) + \mathcal{J}_0(\tilde{\boldsymbol{\theta}}_t - \boldsymbol{\theta}_0) - \mathbf{y}$ denote the residuals.

Assume that the step sizes of the gradient descent ascent updates satisfy $\frac{\mu}{\eta} \geq 8\beta^2$ as well as $0 <$

$\mu \leq 1$. Moreover, assume that the assumptions 1-3 hold for the Jacobian $\mathcal{J}(\boldsymbol{\theta})$ of $f(\boldsymbol{\theta})$ around the initialization $\boldsymbol{\theta}_0 \in \mathbb{R}^n$ with parameters α, β, ϵ , and

$$R := 2\gamma \frac{\beta^2}{\alpha^2} \left\| \begin{bmatrix} \mathcal{J}_0^\dagger & 0 \\ 0 & \mathcal{J}_0^\dagger \end{bmatrix} \mathbf{z}_0 \right\|_{\ell_2} + \frac{18\epsilon\beta^2\gamma^2}{\alpha^4} \|\mathbf{z}_0\|_{\ell_2}, \quad (24)$$

which satisfy $4\gamma\beta\epsilon \leq \alpha^2$. Here, $1 \leq \gamma \leq 5$ is a constant, which only depends on μ, η , and \mathcal{J}_0 . By \mathcal{J}_0^\dagger we denote the pseudo-inverse of the Jacobian at initialization \mathcal{J}_0 . Then, assuming the same initialization $\boldsymbol{\theta}_0 = \tilde{\boldsymbol{\theta}}_0, \mathbf{d}_0 = \tilde{\mathbf{d}}_0$ (and, hence, $\mathbf{z}_0 = \tilde{\mathbf{z}}_0$), the following holds for all iterations $t \geq 0$.

- $\|\mathbf{z}_t\|_{\ell_2}$ converges to 0 with a geometric rate, i.e.

$$\|\mathbf{z}_t\|_{\ell_2} \leq \gamma \left(1 - \frac{\eta\alpha^2}{2}\right)^t \|\mathbf{z}_0\|_{\ell_2}. \quad (25)$$

- The trajectories of \mathbf{z}_t and $\tilde{\mathbf{z}}_t$ stay close to each other and converge to the same limit, i.e.

$$\begin{aligned} \|\mathbf{z}_t - \tilde{\mathbf{z}}_t\|_{\ell_2} &\leq 2\eta\gamma^2\beta\epsilon \cdot t \left(1 - \frac{\eta\alpha^2}{2}\right)^{t-1} \|\mathbf{z}_0\|_{\ell_2} \\ &\leq \frac{4\gamma^2\beta\epsilon}{e \left(15 \ln \frac{16}{15}\right) \alpha^2} \|\mathbf{z}_0\|_{\ell_2}. \end{aligned} \quad (26)$$

- The parameters of the original and linearized problems stay close to each other, i.e.

$$\|\tilde{\boldsymbol{\theta}}_t - \boldsymbol{\theta}_t\|_{\ell_2} \leq \frac{9\epsilon\beta^2\gamma^2}{\alpha^4} \|\mathbf{z}_0\|_{\ell_2}, \quad (27)$$

- The parameters of the original problem stay close to the initialization, i.e.

$$\|\boldsymbol{\theta}_t - \boldsymbol{\theta}_0\|_{\ell_2} \leq \frac{R}{2}. \quad (28)$$

Theorem A.4 will be the main ingredient in the proof of Theorem 2.1. However, as discussed in Section 2.4 we believe that this meta theorem can be used to deal with a much richer class of generators and discriminators.

A.3.1 PROOF OF THEOREM A.4

We will prove the statements in the theorem by induction. The base case for $\tau = 0$ is trivial. Now assume that the equations equation 25 to equation 28 hold for $\tau = 0, \dots, t-1$. Our goal is to show that they hold for iteration t as well.

Part I: First, we are going to show that $\boldsymbol{\theta}_t \in \mathcal{B}_R(\boldsymbol{\theta}_0)$. Note that by the triangle inequality and the induction assumption we have that

$$\begin{aligned} \|\boldsymbol{\theta}_t - \boldsymbol{\theta}_0\|_{\ell_2} &\leq \|\boldsymbol{\theta}_t - \boldsymbol{\theta}_{t-1}\|_{\ell_2} + \|\boldsymbol{\theta}_{t-1} - \boldsymbol{\theta}_0\|_{\ell_2} \\ &\leq \|\boldsymbol{\theta}_t - \boldsymbol{\theta}_{t-1}\|_{\ell_2} + \frac{R}{2}. \end{aligned}$$

Hence, in order to prove the claim it remains to show that $\|\boldsymbol{\theta}_t - \boldsymbol{\theta}_{t-1}\|_{\ell_2} \leq \frac{R}{2}$. For that, we compute

$$\begin{aligned} \frac{1}{\eta} \|\boldsymbol{\theta}_t - \boldsymbol{\theta}_{t-1}\|_{\ell_2} &= \|\mathcal{J}^T(\boldsymbol{\theta}_{t-1}) \mathbf{d}_{t-1}\|_{\ell_2} \\ &\leq \|\mathcal{J}^T(\boldsymbol{\theta}_{t-1}) \tilde{\mathbf{d}}_{t-1}\|_{\ell_2} + \|\mathcal{J}^T(\boldsymbol{\theta}_{t-1})\| \|\mathbf{d}_{t-1} - \tilde{\mathbf{d}}_{t-1}\|_{\ell_2} \\ &\leq \|\mathcal{J}_0^T \tilde{\mathbf{d}}_{t-1}\|_{\ell_2} + \|\mathcal{J}(\boldsymbol{\theta}_{t-1}) - \mathcal{J}_0\| \|\tilde{\mathbf{d}}_{t-1}\|_{\ell_2} + \|\mathcal{J}^T(\boldsymbol{\theta}_{t-1})\| \|\mathbf{d}_{t-1} - \tilde{\mathbf{d}}_{t-1}\|_{\ell_2} \\ &\stackrel{(i)}{\leq} \gamma \left\| \begin{bmatrix} \mathcal{J}_0^T & 0 \\ 0 & \mathcal{J}_0^T \end{bmatrix} \mathbf{z}_0 \right\|_{\ell_2} + \epsilon \cdot \gamma \|\mathbf{z}_0\|_{\ell_2} + \frac{4\beta^2\epsilon\gamma^2}{e \left(15 \ln \frac{16}{15}\right) \alpha^2} \|\mathbf{z}_0\|_{\ell_2} \\ &\stackrel{(ii)}{\leq} \gamma\beta^2 \left\| \begin{bmatrix} \mathcal{J}_0^\dagger & 0 \\ 0 & \mathcal{J}_0^\dagger \end{bmatrix} \mathbf{z}_0 \right\|_{\ell_2} + \frac{3\beta^2\epsilon\gamma^2}{\alpha^2} \|\mathbf{z}_0\|_{\ell_2}, \end{aligned}$$

where $\gamma \leq 5$ is a constant. Let us verify the last two inequalities. Inequality (ii) holds because $1 \leq \gamma$, $1 \leq \frac{\beta^2}{\alpha^2}$, and

$$\begin{aligned} \left\| \begin{bmatrix} \mathcal{J}_0^T & 0 \\ 0 & \mathcal{J}_0^T \end{bmatrix} \mathbf{z}_0 \right\|_{\ell_2} &= \left\| \begin{bmatrix} \mathbf{V}\Sigma^T\mathbf{W}^T & 0 \\ 0 & \mathbf{V}\Sigma^T\mathbf{W}^T \end{bmatrix} \mathbf{z}_0 \right\|_{\ell_2} \\ &= \sqrt{\sum_{i=1}^n \sigma_i^2 (\langle \mathbf{w}_i, \mathbf{r}_0 \rangle^2 + \langle \mathbf{w}_i, \mathbf{d}_0 \rangle^2)} \\ &\leq \beta^2 \sqrt{\sum_{i=1}^n \frac{1}{\sigma_i^2} (\langle \mathbf{w}_i, \mathbf{r}_0 \rangle^2 + \langle \mathbf{w}_i, \mathbf{d}_0 \rangle^2)} = \beta^2 \left\| \begin{bmatrix} \mathcal{J}_0^\dagger & 0 \\ 0 & \mathcal{J}_0^\dagger \end{bmatrix} \mathbf{z}_0 \right\|_{\ell_2}. \quad (29) \end{aligned}$$

Also (i) follows from assumptions 1-3, $\|\mathbf{d}_{t-1} - \tilde{\mathbf{d}}_{t-1}\|_{\ell_2} \leq \|\mathbf{z}_{t-1} - \tilde{\mathbf{z}}_{t-1}\|_{\ell_2}$ together with induction assumption equation 26, $\|\tilde{\mathbf{d}}_{t-1}\|_{\ell_2} \leq \|\tilde{\mathbf{z}}_{t-1}\|_{\ell_2} \leq \|\mathbf{z}_0\|_{\ell_2}$, and

$$\begin{aligned} \left\| \mathcal{J}_0^T \tilde{\mathbf{d}}_{t-1} \right\|_{\ell_2} &\leq \left\| \begin{bmatrix} \mathcal{J}_0^T \tilde{\mathbf{r}}_{t-1} \\ \mathcal{J}_0^T \tilde{\mathbf{d}}_{t-1} \end{bmatrix} \right\|_{\ell_2} \\ &= \left\| \begin{bmatrix} \mathbf{I} & -\eta \mathcal{J}_0^T \mathcal{J}_0 \\ \mu \mathbf{I} & (1-\mu) \mathbf{I} \end{bmatrix} \begin{bmatrix} \mathcal{J}_0^T \tilde{\mathbf{r}}_{t-2} \\ \mathcal{J}_0^T \tilde{\mathbf{d}}_{t-2} \end{bmatrix} \right\|_{\ell_2} \\ &= \left\| \begin{bmatrix} \mathbf{I} & -\eta \mathcal{J}_0^T \mathcal{J}_0 \\ \mu \mathbf{I} & (1-\mu) \mathbf{I} \end{bmatrix}^{t-1} \begin{bmatrix} \mathcal{J}_0^T \tilde{\mathbf{r}}_0 \\ \mathcal{J}_0^T \tilde{\mathbf{d}}_0 \end{bmatrix} \right\|_{\ell_2} \leq \gamma (1 - \eta \alpha^2)^{t-1} \left\| \begin{bmatrix} \mathcal{J}_0^T & 0 \\ 0 & \mathcal{J}_0^T \end{bmatrix} \mathbf{z}_0 \right\|_{\ell_2}, \quad (30) \end{aligned}$$

where in the last inequality we applied Lemma A.3. Finally, by using $\eta \leq \frac{1}{8\beta^2}$ we arrive at

$$\begin{aligned} \|\boldsymbol{\theta}_t - \boldsymbol{\theta}_{t-1}\|_{\ell_2} &\leq \gamma \eta \beta^2 \left\| \begin{bmatrix} \mathcal{J}_0^\dagger & 0 \\ 0 & \mathcal{J}_0^\dagger \end{bmatrix} \mathbf{z}_0 \right\|_{\ell_2} + \frac{3\eta \beta^2 \epsilon \gamma^2}{\alpha^2} \|\mathbf{z}_0\|_{\ell_2} \\ &\leq \frac{\gamma}{8} \left\| \begin{bmatrix} \mathcal{J}_0^\dagger & 0 \\ 0 & \mathcal{J}_0^\dagger \end{bmatrix} \mathbf{z}_0 \right\|_{\ell_2} + \frac{3\epsilon \gamma^2}{8\alpha^2} \|\mathbf{z}_0\|_{\ell_2} \\ &\leq \frac{R}{2}, \end{aligned}$$

where the last line is directly due to inequality (24), $\gamma \leq 5$, and $\alpha \leq \beta$. Hence, we have established $\boldsymbol{\theta}_t \in \mathcal{B}_R(\boldsymbol{\theta}_0)$.

Part II: In Lemma A.3 we showed that the time invariant system of state equations $\tilde{\mathbf{z}}_{t+1} = \mathbf{A}\tilde{\mathbf{z}}_t$ is uniformly exponentially stable, i.e. $\|\tilde{\mathbf{z}}_t\|_{\ell_2}$ goes down to zero exponentially fast. Now by using the assumption that the Jacobian remains close to the Jacobian at the initialization \mathcal{J}_0 , we aim to show the exponential stability of the time variant system of the state equations 16. For that, we compute

$$\begin{aligned} \mathbf{z}_t = \mathbf{A}_{t-1} \mathbf{z}_{t-1} &= \begin{bmatrix} \mathbf{I} & -\eta \mathcal{J}_{t,t-1} \mathcal{J}_{t-1}^T \\ \mu \mathbf{I} & (1-\mu) \mathbf{I} \end{bmatrix} \mathbf{z}_{t-1} \\ &= \begin{bmatrix} \mathbf{I} & -\eta \mathcal{J}_0 \mathcal{J}_0^T \\ \mu \mathbf{I} & (1-\mu) \mathbf{I} \end{bmatrix} \mathbf{z}_{t-1} + \begin{bmatrix} \eta (\mathcal{J}_0 \mathcal{J}_0^T - \mathcal{J}_{t,t-1} \mathcal{J}_{t-1}^T) \mathbf{d}_{t-1} \\ \mathbf{0} \end{bmatrix} \\ &=: \mathbf{A} \mathbf{z}_{t-1} + \boldsymbol{\Delta}_{t-1}. \end{aligned}$$

Now set $\lambda := 1 - \eta\alpha^2$. By induction, we obtain the relation $\mathbf{z}_t = \mathbf{A}^t \mathbf{z}_0 + \sum_{i=0}^{t-1} \mathbf{A}^{t-1-i} \Delta_i$. Hence,

$$\begin{aligned}
\|\mathbf{z}_t\|_{\ell_2} &= \left\| \mathbf{A}^t \mathbf{z}_0 + \sum_{i=0}^{t-1} \mathbf{A}^{t-1-i} \Delta_i \right\|_{\ell_2} \\
&\leq \|\mathbf{A}^t \mathbf{z}_0\|_{\ell_2} + \sum_{i=0}^{t-1} \|\mathbf{A}^{t-1-i} \Delta_i\|_{\ell_2} \\
&\leq \gamma \lambda^t \|\mathbf{z}_0\|_{\ell_2} + \sum_{i=0}^{t-1} \gamma \lambda^{t-1-i} \|\eta (\mathcal{J}_0 \mathcal{J}_0^T - \mathcal{J}_{i+1,i} \mathcal{J}_i^T)\| \|\mathbf{d}_i\|_{\ell_2} \\
&\leq \gamma \lambda^t \|\mathbf{z}_0\|_{\ell_2} + \sum_{i=0}^{t-1} \eta \gamma \lambda^{t-1-i} (2\beta\epsilon) \|\mathbf{z}_i\|_{\ell_2}. \tag{31}
\end{aligned}$$

The second inequality holds because of Lemma A.3. The last inequality holds because by combining our assumptions 1 to 3 with $\boldsymbol{\theta}_t \in \mathcal{B}_R(\boldsymbol{\theta}_0)$ and the induction assumption 28 for $0 \leq i \leq t-1$, we have that

$$\begin{aligned}
\|\mathcal{J}_0 \mathcal{J}_0^T - \mathcal{J}_{i+1,i} \mathcal{J}_i^T\| &= \|\mathcal{J}_0 \mathcal{J}_0^T - \mathcal{J}_0 \mathcal{J}_i^T + \mathcal{J}_0 \mathcal{J}_i^T - \mathcal{J}_{i+1,i} \mathcal{J}_i^T\| \\
&\leq \|\mathcal{J}_0\| \|\mathcal{J}_0 - \mathcal{J}_i\| + \|\mathcal{J}_0 - \mathcal{J}_{i+1,i}\| \|\mathcal{J}_i\| \\
&\leq \beta \|\mathcal{J}_0 - \mathcal{J}_i\| + \beta \|\mathcal{J}_0 - \mathcal{J}_{i+1,i}\| \\
&\leq 2\beta\epsilon. \tag{32}
\end{aligned}$$

In order to deal with inequality 31, we will rely on the following lemma.

Lemma A.5 (Rugh, 1996, Lemma 24.5) Consider two real sequences $p(t)$ and $\phi(t)$, where $p(t) \geq 0$ for all $t \geq 0$ and

$$\phi(t) \leq \begin{cases} \psi, & \text{if } t = 0 \\ \psi + \eta \sum_{i=0}^{t-1} p(i) \phi(i), & \text{if } t \geq 1 \end{cases}$$

where η and ψ are constants with $\eta \geq 0$. Then for all $t \geq 1$ we have

$$\phi(t) \leq \psi \prod_{i=0}^{t-1} (1 + \eta \cdot p(i)).$$

Now we define $\phi_t = \frac{\|\mathbf{z}_t\|_{\ell_2}}{\lambda^t}$ and rewrite inequality 31 as

$$\phi_t \leq \gamma \phi_0 + \sum_{i=0}^{t-1} \frac{2\eta\gamma\beta\epsilon}{\lambda} \phi_i.$$

Hence, Lemma A.5 yields that

$$\begin{aligned}
\phi_t &\leq \gamma \phi_0 \prod_{i=0}^{t-1} \left(1 + \frac{2\eta\gamma\beta\epsilon}{\lambda} \right) \\
&= \gamma \phi_0 \left(1 + \frac{2\eta\gamma\beta\epsilon}{\lambda} \right)^t \\
&\stackrel{(i)}{\leq} \gamma \phi_0 \left(1 + \frac{\eta\alpha^2}{2\lambda} \right)^t \\
&\stackrel{(ii)}{=} \gamma \phi_0 \left(\frac{1 - \frac{\eta\alpha^2}{2}}{1 - \eta\alpha^2} \right)^t,
\end{aligned}$$

where (i) follows from $4\gamma\beta\epsilon \leq \alpha^2$ and (ii) holds by inserting $\lambda = 1 - \eta\alpha^2$. Inserting the definition of ϕ_0 and ϕ_t we obtain that

$$\|\mathbf{z}_t\|_{\ell_2} \leq \gamma \left(1 - \frac{\eta\alpha^2}{2}\right)^t \|\mathbf{z}_0\|_{\ell_2}.$$

This completes the proof of Part II.

Part III: In this part, our aim is to show that the error vector $\mathbf{e}_t := \mathbf{z}_t - \tilde{\mathbf{z}}_t$ obeys inequality 26. First, note that

$$\begin{aligned} \mathbf{e}_t &= \mathbf{z}_t - \tilde{\mathbf{z}}_t \\ &\stackrel{(*)}{=} (\mathbf{A}\mathbf{z}_{t-1} + \mathbf{\Delta}_{t-1}) - \mathbf{A}\tilde{\mathbf{z}}_{t-1} \\ &= \mathbf{A}\mathbf{e}_{t-1} + \mathbf{\Delta}_{t-1}, \end{aligned}$$

where in (*) we used the same notation as in Part II for $\mathbf{\Delta}_{t-1}$. Using a recursive argument as well as $\mathbf{e}_0 = 0$ we obtain that

$$\begin{aligned} \|\mathbf{e}_t\|_{\ell_2} &= \left\| \sum_{i=0}^{t-1} \mathbf{A}^{t-1-i} \mathbf{\Delta}_i \right\|_{\ell_2} \\ &\leq \sum_{i=0}^{t-1} \eta\gamma (1 - \eta\alpha^2)^{t-1-i} \|\mathbf{\Delta}_i\|_{\ell_2} \\ &\stackrel{(i)}{\leq} \sum_{i=0}^{t-1} \eta\gamma (1 - \eta\alpha^2)^{t-1-i} \|\mathcal{J}_0 \mathcal{J}_0^T - \mathcal{J}_{i+1,i} \mathcal{J}_i^T\| \|\mathbf{d}_i\|_{\ell_2} \\ &\stackrel{(ii)}{\leq} \sum_{i=0}^{t-1} 2\eta\beta\epsilon\gamma (1 - \eta\alpha^2)^{t-1-i} \|\mathbf{z}_i\|_{\ell_2}. \end{aligned}$$

The first inequality follows from the triangle inequality and Lemma A.3. Inequality (i) follows from the definition of $\mathbf{\Delta}_i$. Inequality (ii) follows from inequality 32. Setting $c := 2\eta\beta\epsilon$ we continue

$$\begin{aligned} \|\mathbf{e}_t\|_{\ell_2} &\leq \sum_{i=0}^{t-1} c\gamma (1 - \eta\alpha^2)^{t-i-1} \|\mathbf{z}_i\|_{\ell_2} \\ &\stackrel{(iii)}{\leq} \sum_{i=0}^{t-1} c\gamma^2 (1 - \eta\alpha^2)^{t-i-1} \left(1 - \frac{\eta\alpha^2}{2}\right)^i \|\mathbf{z}_0\|_{\ell_2} \\ &\stackrel{(iv)}{\leq} \sum_{i=0}^{t-1} c\gamma^2 \left(1 - \frac{\eta\alpha^2}{2}\right)^{t-1} \|\mathbf{z}_0\|_{\ell_2} \\ &= 2\eta\gamma^2\beta\epsilon \cdot t \left(1 - \frac{\eta\alpha^2}{2}\right)^{t-1} \|\mathbf{z}_0\|_{\ell_2}. \end{aligned}$$

Here (iii) holds because of our induction hypothesis 25 and (iv) follows simply from $1 - \eta\alpha^2 \leq 1 - \frac{\eta\alpha^2}{2}$. This shows the first part of equation 26 for iteration t . Finally, to derive the second part of equation 26 we observe that for all $t \geq 0$ and $0 < x \leq \frac{1}{16}$ we have $t(1-x)^{t-1} \leq \frac{1}{e(15 \ln \frac{16}{15})x}$. Since $\frac{\eta\alpha^2}{2} \leq \frac{\mu\alpha^2}{16\beta^2} \leq \frac{1}{16}$ we can use this estimate, which yields

$$\begin{aligned} \|\mathbf{e}_t\|_{\ell_2} &\leq 2\eta\gamma^2\beta\epsilon \cdot t \left(1 - \frac{\eta\alpha^2}{2}\right)^{t-1} \|\mathbf{z}_0\|_{\ell_2} \\ &\leq \frac{4\gamma^2\beta\epsilon}{e(15 \ln \frac{16}{15})\alpha^2} \|\mathbf{z}_0\|_{\ell_2}. \end{aligned}$$

Hence, we have shown equation 26.

Part IV: In this part, we aim to show that the parameters of the original and linearized problems are close. For that, we compute that

$$\begin{aligned}
\frac{1}{\eta} \|\boldsymbol{\theta}_t - \tilde{\boldsymbol{\theta}}_t\|_{\ell_2} &= \left\| \sum_{i=0}^{t-1} \nabla_{\boldsymbol{\theta}} h(\boldsymbol{\theta}_i, \mathbf{d}_i) - \nabla_{\boldsymbol{\theta}} h_{\text{lin}}(\boldsymbol{\theta}_i, \mathbf{d}_i) \right\|_{\ell_2} \\
&= \left\| \sum_{i=0}^{t-1} \mathcal{J}^T(\boldsymbol{\theta}_i) \mathbf{d}_i - \mathcal{J}_0^T \tilde{\mathbf{d}}_i \right\|_{\ell_2} \\
&\leq \sum_{i=0}^{t-1} \left\| (\mathcal{J}^T(\boldsymbol{\theta}_i) - \mathcal{J}_0^T) \tilde{\mathbf{d}}_i \right\|_{\ell_2} + \sum_{i=0}^{t-1} \left\| \mathcal{J}^T(\boldsymbol{\theta}_i) (\mathbf{d}_i - \tilde{\mathbf{d}}_i) \right\|_{\ell_2} \\
&\stackrel{(i)}{\leq} \sum_{i=0}^{t-1} \epsilon \|\tilde{\mathbf{z}}_i\|_{\ell_2} + \beta \sum_{i=0}^{t-1} \|\mathbf{e}_i\|_{\ell_2} \\
&\stackrel{(ii)}{\leq} \gamma \epsilon \sum_{i=0}^{t-1} (1 - \eta\alpha^2)^i \|\mathbf{z}_0\|_{\ell_2} + 2\eta\gamma^2\beta^2\epsilon \sum_{i=0}^{t-1} i \left(1 - \frac{\eta\alpha^2}{2}\right)^{i-1} \|\mathbf{z}_0\|_{\ell_2}.
\end{aligned}$$

Here (i) follows from assumptions 2 and 3, and (ii) holds because of Lemma A.3 and our induction hypothesis 26. Hence, using the formula $\sum_{i=0}^t ix^i = \frac{x(1+tx^{t+1}-(t+1)x^t)}{(x-1)^2}$ we obtain that

$$\begin{aligned}
\frac{1}{\eta} \|\boldsymbol{\theta}_t - \tilde{\boldsymbol{\theta}}_t\|_{\ell_2} &\leq \gamma \epsilon \|\mathbf{z}_0\|_{\ell_2} \left(\frac{1 - (1 - \eta\alpha^2)^t}{\eta\alpha^2} + 2\eta\beta^2\gamma \frac{1 - t \left(1 - \frac{\eta\alpha^2}{2}\right)^{t-1} + (t-1) \left(1 - \frac{\eta\alpha^2}{2}\right)^t}{\left(\frac{\eta\alpha^2}{2}\right)^2} \right) \\
&\leq \gamma \epsilon \|\mathbf{z}_0\|_{\ell_2} \left(\frac{1}{\eta\alpha^2} + 2\eta\beta^2\gamma \frac{1}{\left(\frac{\eta\alpha^2}{2}\right)^2} \right) \\
&\stackrel{(iii)}{\leq} \gamma \epsilon \|\mathbf{z}_0\|_{\ell_2} \left(\frac{\beta^2\gamma}{\eta\alpha^4} + \frac{8\beta^2\gamma}{\eta\alpha^4} \right) \\
&= \frac{9\epsilon\beta^2\gamma^2}{\eta\alpha^4} \|\mathbf{z}_0\|_{\ell_2},
\end{aligned}$$

where (iii) holds due to $1 \leq \gamma$ and $1 \leq \frac{\beta^2}{\alpha^2}$. Hence, we have established inequality 27 for iteration t .

Part V: In this part, we are going to prove equation 28 for iteration t . First, it follows from the triangle inequality that

$$\begin{aligned}
\|\boldsymbol{\theta}_t - \boldsymbol{\theta}_0\|_{\ell_2} &\leq \|\tilde{\boldsymbol{\theta}}_t - \boldsymbol{\theta}_0\|_{\ell_2} + \|\boldsymbol{\theta}_t - \tilde{\boldsymbol{\theta}}_t\|_{\ell_2} \\
&\leq \|\tilde{\boldsymbol{\theta}}_t - \boldsymbol{\theta}_0\|_{\ell_2} + \frac{9\epsilon\beta^2\gamma^2}{\alpha^4} \|\mathbf{z}_0\|_{\ell_2},
\end{aligned}$$

where in the second inequality we have used Part IV. Now we bound $\|\tilde{\boldsymbol{\theta}}_t - \boldsymbol{\theta}_0\|_{\ell_2}$ from above as follows

$$\begin{aligned}
\|\tilde{\boldsymbol{\theta}}_t - \boldsymbol{\theta}_0\|_{\ell_2} &= \eta \left\| \sum_{i=0}^{t-1} \mathcal{J}_0^T \tilde{\mathbf{d}}_i \right\|_{\ell_2} \\
&\leq \eta \sum_{i=0}^{t-1} \|\mathcal{J}_0^T \tilde{\mathbf{d}}_i\|_{\ell_2} \\
&\stackrel{(i)}{\leq} \eta \gamma \sum_{i=0}^{t-1} (1 - \eta \alpha^2)^i \left\| \begin{bmatrix} \mathcal{J}_0^T & 0 \\ 0 & \mathcal{J}_0^T \end{bmatrix} \mathbf{z}_0 \right\|_{\ell_2} \\
&= \eta \gamma \frac{1 - (1 - \eta \alpha^2)^t}{\eta \alpha^2} \left\| \begin{bmatrix} \mathcal{J}_0^T & 0 \\ 0 & \mathcal{J}_0^T \end{bmatrix} \mathbf{z}_0 \right\|_{\ell_2} \\
&\stackrel{(ii)}{\leq} \gamma \frac{\beta^2}{\alpha^2} \left\| \begin{bmatrix} \mathcal{J}_0^\dagger & 0 \\ 0 & \mathcal{J}_0^\dagger \end{bmatrix} \mathbf{z}_0 \right\|_{\ell_2},
\end{aligned}$$

where (i) holds by equation 30 and (ii) holds by equation 29. Hence, it follows from the definition of R (equation 24) that

$$\begin{aligned}
\|\boldsymbol{\theta}_t - \boldsymbol{\theta}_0\|_{\ell_2} &\leq \gamma \frac{\beta^2}{\alpha^2} \left\| \begin{bmatrix} \mathcal{J}_0^\dagger & 0 \\ 0 & \mathcal{J}_0^\dagger \end{bmatrix} \mathbf{z}_0 \right\|_{\ell_2} + \frac{9\epsilon\beta^2\gamma^2}{\alpha^4} \|\mathbf{z}_0\|_{\ell_2} \\
&= \frac{R}{2}.
\end{aligned}$$

This completes the proof.

A.4 PRELIMINARIES FOR PROOFS OF RESULTS WITH ONE-HIDDEN LAYER GENERATOR AND LINEAR DISCRIMINATOR

In this section, we gather some preliminary results that will be useful in proving the main results i.e. Theorems 2.1 and 2.2. We begin by noting that Theorem 2.1 is an instance of Theorem A.4 with $f(\mathbf{W}) = \frac{1}{n} \sum_{i=1}^n \mathbf{V} \cdot \phi(\mathbf{W} \mathbf{z}_i)$. We thus begin this section by noting that $f(\mathbf{W})$ can be rewritten as follows

$$f(\mathbf{W}) = \mathbf{V} \cdot \begin{bmatrix} \frac{1}{n} \sum_{i=1}^n \phi(\mathbf{w}_1^T \mathbf{z}_i) \\ \vdots \\ \frac{1}{n} \sum_{i=1}^n \phi(\mathbf{w}_k^T \mathbf{z}_i) \end{bmatrix}.$$

Furthermore, the Jacobian of this mapping $f(\mathbf{W})$ takes the form

$$\mathcal{J}(\mathbf{W}) = \frac{1}{n} \sum_{i=1}^n (\mathbf{V} \cdot \text{diag}(\phi'(\mathbf{W} \mathbf{z}_i))) \otimes \mathbf{z}_i^T. \quad (33)$$

To characterize the spectral properties of this Jacobian it will be convenient to write down the expression for $\mathcal{J}(\mathbf{W})\mathcal{J}(\mathbf{W})^T$ which has a compact form

$$\begin{aligned}
\mathcal{J}(\mathbf{W})\mathcal{J}(\mathbf{W})^T &\stackrel{(i)}{=} \frac{1}{n^2} \sum_{i,j=1}^n \left((\mathbf{V} \cdot \text{diag}(\phi'(\mathbf{W} \mathbf{z}_i))) \otimes \mathbf{z}_i^T \right) \left(\text{diag}(\phi'(\mathbf{W} \mathbf{z}_j)) \mathbf{V}^T \otimes \mathbf{z}_j \right) \\
&\stackrel{(ii)}{=} \frac{1}{n^2} \sum_{i,j=1}^n \left(\mathbf{V} \text{diag}(\phi'(\mathbf{W} \mathbf{z}_i)) \text{diag}(\phi'(\mathbf{W} \mathbf{z}_j)) \mathbf{V}^T \right) \otimes \left(\mathbf{z}_i^T \mathbf{z}_j \right) \\
&= \frac{1}{n^2} \mathbf{V} \text{diag}_{\ell=1, \dots, k} \left(\left\| \sum_{i=1}^n \mathbf{z}_i \phi'(\mathbf{w}_\ell^T \mathbf{z}_i) \right\|_{\ell_2}^2 \right) \mathbf{V}^T \\
&= \frac{1}{n^2} \mathbf{V} \cdot \mathbf{D}^2 \cdot \mathbf{V}^T,
\end{aligned}$$

where D is a diagonal matrix with entries

$$D_{\ell\ell} = \left\| \sum_{i=1}^n z_i \phi'(\mathbf{w}_\ell^T z_i) \right\|_{\ell_2} = \|\mathbf{Z}^T \phi'(\mathbf{Z}\mathbf{w}_\ell)\|_{\ell_2}, \quad (34)$$

and $\mathbf{Z} \in \mathbb{R}^{n \times d}$ contains the z_i 's in its rows. Note that we used simple properties of kronecker product in (i) and (ii), namely $(A \otimes B)^T = A^T \otimes B^T$ and $(A \otimes B)(C \otimes D) = (AC) \otimes (BD)$. The next lemma establishes concentration of the diagonal entries of matrix D^2 around their mean, which will be used in the future lemmas regarding the spectrum of the Jacobian mapping. The proof is deferred to Appendix B.1.

Lemma A.6 *Suppose $\mathbf{w} \in \mathbb{R}^d$ is a fixed vector, $z_1, z_2, \dots, z_n \in \mathbb{R}^d$ are distributed as $\mathcal{N}(0, \sigma_z^2 \mathbf{I}_d)$ and constitute the rows of $\mathbf{Z} \in \mathbb{R}^{n \times d}$. Then for any $0 \leq \delta \leq \frac{3}{2}$ the random variable $D = \|\mathbf{Z}^T \phi'(\mathbf{Z}\mathbf{w})\|_{\ell_2}$ satisfies*

$$(1 - \delta) \mathbb{E}(D^2) \leq D^2 \leq (1 + \delta) \mathbb{E}(D^2)$$

with probability at least $1 - 2 \left(e^{-\frac{n\delta^2}{18}} + e^{-\frac{d\delta^2}{54}} + e^{-c_1 n\delta} \right)$ where c_1 is a positive constant. Moreover we have

$$\mathbb{E}(D^2) = \sigma_z^2 \left(\frac{nd}{2} + \frac{n(n-1)}{2\pi} \right).$$

Furthermore, using the above equation we have

$$\mathbb{E}[\mathcal{J}(\mathbf{W})\mathcal{J}(\mathbf{W})^T] = \frac{\sigma_z^2 \left(d + \frac{n-1}{\pi} \right)}{2n} \mathbf{V}\mathbf{V}^T.$$

A.5 LEMMAS REGARDING THE INITIAL MISFIT AND THE SPECTRUM OF THE JACOBIAN

In this section, we state some lemmas regarding the spectrum of the Jacobian mapping and the initial misfit, and defer their proofs to Appendix B. First, we state a result on the minimum singular value of the Jacobian mapping at initialization.

Lemma A.7 (Minimum singular value of the Jacobian at initialization) *Consider our GAN model with a linear discriminator and the generator being a one hidden layer neural network of the form $\mathbf{z} \mapsto \mathbf{V}\phi(\mathbf{W}\mathbf{z})$, where we have n independent data points $z_1, z_2, \dots, z_n \in \mathbb{R}^d$ distributed as $\mathcal{N}(0, \sigma_z^2 \mathbf{I}_d)$ and aggregated as the rows of a matrix $\mathbf{Z} \in \mathbb{R}^{n \times d}$, and $\mathbf{V} \in \mathbb{R}^{m \times k}$ has i.i.d $\mathcal{N}(0, \sigma_v^2)$ entries. We also assume that $\mathbf{W}_0 \in \mathbb{R}^{k \times d}$ has i.i.d $\mathcal{N}(0, \sigma_w^2)$ entries and all entries of \mathbf{W}_0 , \mathbf{V} , and \mathbf{Z} are independent. Then the Jacobian matrix at the initialization point obeys*

$$\sigma_{\min}(\mathcal{J}(\mathbf{W}_0)) \geq \left(\sqrt{(1-\delta)^2 k - (1+\delta)^2} - \sqrt{m} (1+\eta) (1+\delta) \right) \sigma_v \sigma_z \sqrt{\frac{d + \frac{n-1}{\pi}}{2n}}, \quad 0 \leq \delta \leq \frac{3}{2}$$

with probability at least $1 - 3e^{-\frac{\eta^2 m}{8}} - 2k \cdot \left(e^{-\frac{n\delta^2}{18}} + e^{-\frac{d\delta^2}{54}} + e^{-c_1 n\delta} \right)$, where c_1 is a positive constant.

Next lemma helps us bound the spectral norm of the Jacobian at initialization, which will be used later to derive upper bounds on Jacobian at every point near initialization.

Lemma A.8 (spectral norm of the Jacobian at initialization) *Following the setup of previous lemma, the operator norm of the Jacobian matrix at initialization point $\mathbf{W}_0 \in \mathbb{R}^{k \times d}$ satisfies*

$$\|\mathcal{J}(\mathbf{W}_0)\| \leq (1 + \delta) \sigma_v \sigma_z \left(\sqrt{k} + 2\sqrt{m} \right) \sqrt{\frac{d + \frac{n-1}{\pi}}{2n}}, \quad 0 \leq \delta \leq \frac{3}{2}$$

with probability at least $1 - e^{-\frac{m}{2}} - k \cdot \left(e^{-\frac{n\delta^2}{18}} + e^{-\frac{d\delta^2}{54}} + e^{-c_1 n\delta} \right)$, with c_1 a positive constant.

The next lemma is adapted from [Van Veen et al. \(2018\)](#) and allows us to bound the variations in the Jacobian matrix around initialization.

Lemma A.9 (single-sample Jacobian perturbation) Let $\mathbf{V} \in \mathbb{R}^{m \times k}$ be a matrix with i.i.d. $\mathcal{N}(0, \sigma_v^2)$ entries, $\mathbf{W} \in \mathbb{R}^{k \times d}$, and define the Jacobian mapping $\mathcal{J}(\mathbf{W}; z) = (\mathbf{V} \text{diag}(\phi'(\mathbf{W}z))) \otimes z^T$. Then, by taking \mathbf{W}_0 to be a random matrix with i.i.d. $\mathcal{N}(0, \sigma_w^2)$ entries, we have

$$\|\mathcal{J}(\mathbf{W}; z) - \mathcal{J}(\mathbf{W}_0; z)\| \leq \sigma_v \|z\|_{\ell_2} \left(2\sqrt{m} + \sqrt{6 \left(\frac{2kR}{\sigma_w} \right)^{\frac{2}{3}} \log \left(\frac{k}{3 \left(\frac{2kR}{\sigma_w} \right)^{\frac{2}{3}}} \right)} \right)$$

for all $\mathbf{W} \in \mathbb{R}^{k \times d}$ obeying $\|\mathbf{W} - \mathbf{W}_0\| \leq R$ with probability at least $1 - e^{-\frac{m}{2}} - e^{-\frac{(2kR)}{\sigma_w} \frac{2}{3}}$.

Our final key lemma bounds the initial misfit $f(\mathbf{W}_0) - \mathbf{y} := \frac{1}{n} \sum_{i=1}^n \mathbf{V} \phi(\mathbf{W}_0 \mathbf{z}_i) - \bar{\mathbf{x}}$.

Lemma A.10 (Initial misfit) Consider our GAN model with a linear discriminator and the generator being a one hidden layer neural network of the form $\mathbf{z} \mapsto \mathbf{V} \phi(\mathbf{W}z)$, where we have n independent data points $\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_n \in \mathbb{R}^d$ distributed as $\mathcal{N}(0, \sigma_z^2 \mathbf{I}_d)$ and aggregated as the rows of a matrix $\mathbf{Z} \in \mathbb{R}^{n \times d}$, and $\mathbf{V} \in \mathbb{R}^{m \times k}$ has i.i.d. $\mathcal{N}(0, \sigma_v^2)$ entries. We also assume that the initial $\mathbf{W}_0 \in \mathbb{R}^{k \times d}$ has i.i.d. $\mathcal{N}(0, \sigma_w^2)$ entries. Then the following event

$$\left\| \frac{1}{n} \sum_{i=1}^n \mathbf{V} \phi(\mathbf{W}_0 \mathbf{z}_i) - \bar{\mathbf{x}} \right\|_{\ell_2} \leq (1 + \delta) \frac{1}{\sqrt{2\pi}} \sigma_v \sigma_w \sigma_z \sqrt{kdm} + \|\bar{\mathbf{x}}\|_{\ell_2}, \quad 0 \leq \delta \leq 3$$

holds with probability at least $1 - \left(k \cdot e^{-c_2 n (\delta/27)^2} + e^{-\frac{(\delta/9)^2 m}{2}} + e^{-\frac{(\delta/3)^2 kd}{2}} \right)$, with c_2 a fixed constant.

A.6 PROOF OF THEOREM 2.1

In this section, we prove Theorem 2.1 by using our general meta Theorem A.4. To do this we need to check that Assumptions 1-3 are satisfied with high probability. Specifically, in our case the parameter θ is the matrix \mathbf{W} and the non-linear mapping f is given by $f(\mathbf{W}) = \frac{1}{n} \sum_{i=1}^n \mathbf{V} \phi(\mathbf{W} \mathbf{z}_i)$. We note that in our result $\mathbf{d}_0 = \mathbf{0}$ and thus $\|\mathbf{z}_0\|_{\ell_2} = \|\mathbf{r}_0\|_{\ell_2}$, which simplifies our analysis.

To prove Assumption 1 note that by setting $\delta = \frac{1}{2}$ and $\eta = \frac{1}{3}$ in Lemma A.7, we have

$$\begin{aligned} \sigma_{\min}(\mathcal{J}(\mathbf{W}_0)) &\geq \sigma_v \sigma_z \left(\frac{1}{2} \sqrt{k-9} - 2\sqrt{m} \right) \sqrt{\frac{d + \frac{n-1}{\pi}}{2n}} \\ &=: \alpha. \end{aligned}$$

This holds with probability at least $1 - 3e^{-\frac{m}{2}} - 4k \cdot e^{-c \cdot n} - 2k \cdot e^{-\frac{d}{216}}$, concluding the proof of Assumption 1. Next, by setting $\delta = \frac{1}{2}$ in Lemma A.8 we have

$$\|\mathcal{J}(\mathbf{W}_0)\| \leq \zeta := \frac{3}{2} \sigma_v \sigma_z \left(\sqrt{k} + 2\sqrt{m} \right) \sqrt{\frac{d + \frac{n-1}{\pi}}{2n}}$$

with probability at least $1 - e^{-\frac{m}{2}} - 2k \cdot e^{-c \cdot n} - k \cdot e^{-\frac{d}{216}}$. Now to bound spectral norm of Jacobian at \mathbf{W} where $\|\mathbf{W} - \mathbf{W}_0\| \leq R$ (the value of R is defined in the proof of assumption 3 below), we use triangle inequality to get

$$\|\mathcal{J}(\mathbf{W})\| \leq \|\mathcal{J}(\mathbf{W}_0)\| + \|\mathcal{J}(\mathbf{W}) - \mathcal{J}(\mathbf{W}_0)\|.$$

This last inequality together with assumption 3, which we will prove below, yields

$$\|\mathcal{J}(\mathbf{W})\| \leq \|\mathcal{J}(\mathbf{W}_0)\| + \epsilon \leq \|\mathcal{J}(\mathbf{W}_0)\| + \frac{\alpha^2}{4\gamma\beta} \leq \|\mathcal{J}(\mathbf{W}_0)\| + \frac{\|\mathcal{J}(\mathbf{W}_0)\|^2}{4\beta}.$$

Therefore by choosing $\beta = 2\zeta$ we arrive at

$$\begin{aligned}
\|\mathcal{J}(\mathbf{W})\| &\leq \|\mathcal{J}(\mathbf{W}_0)\| + \frac{\|\mathcal{J}(\mathbf{W}_0)\|^2}{4\beta} \\
&= \|\mathcal{J}(\mathbf{W}_0)\| + \frac{\|\mathcal{J}(\mathbf{W}_0)\|^2}{8\zeta} \\
&\leq \|\mathcal{J}(\mathbf{W}_0)\| + \frac{\|\mathcal{J}(\mathbf{W}_0)\|^2}{8\|\mathcal{J}(\mathbf{W}_0)\|} \\
&\leq 2\|\mathcal{J}(\mathbf{W}_0)\| \\
&\leq 2\zeta = \beta,
\end{aligned}$$

establishing that assumption 2 holds with

$$\beta = 3\sigma_v\sigma_z \left(\sqrt{k} + 2\sqrt{m}\right) \sqrt{\frac{d + \frac{n-1}{\pi}}{2n}}$$

with probability at least $1 - e^{-\frac{m}{2}} - 2k \cdot e^{-c \cdot n} - k \cdot e^{-\frac{d}{216}}$.

Finally to show that Assumption 3 holds, we use the single-sample Jacobian perturbation result of Lemma A.9 combined with the triangle inequality to conclude that

$$\begin{aligned}
\|\mathcal{J}(\mathbf{W}) - \mathcal{J}(\mathbf{W}_0)\| &= \left\| \frac{1}{n} \left(\sum_{i=1}^n \mathcal{J}(\mathbf{W}; z_i) - \mathcal{J}(\mathbf{W}_0; z_i) \right) \right\| \\
&\leq \frac{1}{n} \sum_{i=1}^n \|\mathcal{J}(\mathbf{W}; z_i) - \mathcal{J}(\mathbf{W}_0; z_i)\| \\
&\leq \frac{\sigma_v}{n} \left(\sum_{i=1}^n \|z_i\|_{\ell_2} \right) \left(2\sqrt{m} + \sqrt{6 \left(\frac{2kR}{\sigma_w} \right)^{\frac{2}{3}} \log \left(\frac{k}{3 \left(\frac{2kR}{\sigma_w} \right)^{\frac{2}{3}}} \right)} \right) \\
&\stackrel{(i)}{\leq} \sigma_v \frac{\|\mathbf{Z}\|_F}{\sqrt{n}} \left(2\sqrt{m} + \sqrt{6 \left(\frac{2kR}{\sigma_w} \right)^{\frac{2}{3}} \log \left(\frac{k}{3 \left(\frac{2kR}{\sigma_w} \right)^{\frac{2}{3}}} \right)} \right) \\
&\stackrel{(ii)}{\leq} \frac{5}{4} \sigma_v \sigma_z \sqrt{d} \left(2\sqrt{m} + \sqrt{6 \left(\frac{2kR}{\sigma_w} \right)^{\frac{2}{3}} \log \left(\frac{k}{3 \left(\frac{2kR}{\sigma_w} \right)^{\frac{2}{3}}} \right)} \right), \quad (35)
\end{aligned}$$

where (i) holds by Cauchy–Schwarz inequality, and (ii) holds because for a Gaussian matrix $\mathbf{Z} \in \mathbb{R}^{n \times d}$ with $\mathcal{N}(0, \sigma_z^2)$ entries the following holds

$$\mathbb{P} \left(\|\mathbf{Z}\|_F \leq \frac{5}{4} \sigma_z \sqrt{nd} \right) \geq \mathbb{P} \left(\|\mathbf{Z}\|_F^2 \leq \frac{3}{2} \sigma_z^2 nd \right) \geq 1 - e^{-\frac{nd}{24}}.$$

Now we set $\epsilon = \frac{\alpha^2}{4\gamma\beta}$ and show that Assumption 3 holds with this choice of ϵ and with radius \tilde{R} , whose value will be defined later in the proof. First, note that

$$\begin{aligned}\epsilon &= \frac{\alpha^2}{4\gamma\beta} \\ &= \frac{\sigma_v^2 \sigma_z^2 \left(\frac{1}{2}\sqrt{k-9} - 2\sqrt{m}\right)^2 \left(\frac{d+\frac{n-1}{2n}}{2n}\right)}{12\gamma\sigma_v\sigma_z \left(\sqrt{k} + 2\sqrt{m}\right) \sqrt{\frac{d+\frac{n-1}{2n}}{2n}}} \\ &\stackrel{(i)}{\geq} \frac{\sigma_v\sigma_z \left(\frac{1}{8}\sqrt{k}\right)^2 \cdot \sqrt{\frac{1}{4\pi}}}{60 \left(3\sqrt{k}\right)} \\ &\geq \frac{\sigma_v\sigma_z\sqrt{k}}{42000},\end{aligned}$$

where (i) holds by assuming $k \geq C \cdot m$ with C being a large positive constant. Combining the last inequality with equation 35, we observe that a sufficient condition for assumption 3 to hold is

$$\frac{5}{4}\sigma_v\sigma_z\sqrt{d} \left(2\sqrt{m} + \sqrt{6 \left(\frac{2kR}{\sigma_w}\right)^{\frac{2}{3}} \log \left(\frac{k}{3 \left(\frac{2kR}{\sigma_w}\right)^{\frac{2}{3}}}\right)} \right) \leq \frac{\sigma_v\sigma_z\sqrt{k}}{42000},$$

which is equivalent to

$$105000\sqrt{md} + 52500 \cdot \sqrt{d} \cdot \sqrt{6 \left(\frac{2kR}{\sigma_w}\right)^{\frac{2}{3}} \log \left(\frac{k}{3 \left(\frac{2kR}{\sigma_w}\right)^{\frac{2}{3}}}\right)} \leq \sqrt{k}.$$

Now the first term in the L.H.S. is upper bounded by $\frac{1}{2}\sqrt{k}$ if $k \geq (210000)^2 md$, and for the second term we need

$$105000 \cdot \sqrt{d} \cdot \sqrt{6 \left(\frac{2kR}{\sigma_w}\right)^{\frac{2}{3}} \log \left(\frac{k}{3 \left(\frac{2kR}{\sigma_w}\right)^{\frac{2}{3}}}\right)} \leq \sqrt{k},$$

which by defining $x = 3d \left(\frac{2R}{\sigma_w\sqrt{k}}\right)^{\frac{2}{3}}$ is equivalent to

$$x \log \frac{d}{x} \leq \frac{1}{2 \cdot 105000^2}.$$

This last inequality holds for $x \leq \frac{c}{\log d}$ with $c < 1$ a sufficiently small positive constant, which translates into

$$R \leq c \frac{\sigma_w\sqrt{k}}{(d \log d)^{\frac{3}{2}}}. \quad (36)$$

So far we have shown that Assumption 3 holds with $\epsilon = \frac{\alpha^2}{4\gamma\beta}$ and with radius \tilde{R} defined as $\tilde{R} := c \frac{\sigma_w\sqrt{k}}{(d \log d)^{\frac{3}{2}}}$, and we conclude that it holds for any radius R less than \tilde{R} as well. Now we work with

the definition of R in equation 24 to show that $R \leq \tilde{R}$:

$$\begin{aligned}
\frac{R}{2} &= \gamma \frac{\beta^2}{\alpha^2} \left\| \begin{bmatrix} \mathcal{J}_0^\dagger & 0 \\ 0 & \mathcal{J}_0^\dagger \end{bmatrix} \mathbf{z}_0 \right\|_{\ell_2} + \frac{9\epsilon\beta^2\gamma^2}{\alpha^4} \|\mathbf{z}_0\|_{\ell_2} \\
&\stackrel{(i)}{\leq} \gamma \frac{\beta^2}{\alpha^3} \|\mathbf{r}_0\|_{\ell_2} + \frac{9\frac{\alpha^2}{4\gamma\beta}\beta^2\gamma^2}{\alpha^4} \|\mathbf{r}_0\|_{\ell_2} \\
&= \gamma \|\mathbf{r}_0\|_{\ell_2} \left(\frac{\beta^2}{\alpha^3} + \frac{9}{4} \frac{\beta}{\alpha^2} \right) \\
&\stackrel{(ii)}{\leq} 20 \frac{\beta^2}{\alpha^3} \|\mathbf{r}_0\|_{\ell_2} \\
&= 20 \frac{\left(3\sigma_v\sigma_z (\sqrt{k} + 2\sqrt{m}) \sqrt{\frac{d+\frac{n-1}{2n}}{2n}} \right)^2}{\left(\sigma_v\sigma_z (\frac{1}{2}\sqrt{k} - 9 - 2\sqrt{m}) \sqrt{\frac{d+\frac{n-1}{2n}}{2n}} \right)^3} \|\mathbf{r}_0\|_{\ell_2} \\
&\stackrel{(iii)}{\leq} C \cdot \frac{1}{\sigma_v\sigma_z\sqrt{k}} \cdot \left(\frac{2}{3}\sigma_v\sigma_w\sigma_z\sqrt{k \cdot d \cdot m} + \|\bar{\mathbf{x}}\|_{\ell_2} \right)
\end{aligned}$$

where (i) holds because $\|\mathcal{J}_0^\dagger\| \leq \frac{1}{\alpha}$ and $4\gamma\beta\epsilon = \alpha^2$, (ii) holds as $1 \leq \frac{\beta}{\alpha}$ and as we substitute $\gamma = 5$ from Lemma A.3, and (iii) follows from $k \geq C \cdot m$ and using $\delta = \frac{1}{3}$ in Lemma A.10. Now a sufficient condition for equation 36 to hold is that

$$\frac{1}{\sigma_v\sigma_z\sqrt{k}} \cdot \left(\frac{2}{3}\sigma_v\sigma_w\sigma_z\sqrt{k \cdot d \cdot m} + \|\bar{\mathbf{x}}\|_{\ell_2} \right) \leq c \frac{\sigma_w\sqrt{k}}{(d \log d)^{\frac{3}{2}}},$$

which is equivalent to

$$\frac{2}{3}\sigma_v\sigma_w\sigma_z \cdot (d \log d)^{\frac{3}{2}} \sqrt{k \cdot d \cdot m} + (d \log d)^{\frac{3}{2}} \|\bar{\mathbf{x}}\|_{\ell_2} \leq c \cdot k\sigma_v\sigma_w\sigma_z.$$

Finally, this inequality is satisfied by assuming $k \geq C \cdot md^4 \log(d)^3$ and setting $\sigma_v\sigma_w\sigma_z \geq \frac{\|\bar{\mathbf{x}}\|_{\ell_2}}{md^{\frac{3}{2}} \log d^{\frac{3}{2}}}$. This shows that assumption 3 holds with probability at least $1 - ne^{-\frac{m}{2}} - ne^{-c \cdot md^3 \log(d)^2} - k \cdot e^{-c \cdot n} - e^{-\frac{m}{1500}} - e^{-\frac{kd}{162}}$, concluding the proof of Theorem 2.1.

A.7 PROOF OF THEOREM 2.2

Consider a nonlinear least-squares optimization problem of the form

$$\min_{\theta \in \mathbb{R}^p} \mathcal{L}(\theta) := \frac{1}{2} \|f(\theta) - y\|_{\ell_2}^2$$

with $f : \mathbb{R}^p \mapsto \mathbb{R}^m$ and $y \in \mathbb{R}^m$. Suppose the Jacobian mapping associated with f satisfies the following three assumptions.

Assumption 1 We assume $\sigma_{\min}(\mathcal{J}(\theta_0)) \geq 2\alpha$ for a fixed point $\theta_0 \in \mathbb{R}^p$.

Assumption 2 Let $\|\cdot\|$ be a norm dominated by the Frobenius norm i.e. $\|\theta\| \leq \|\theta\|_F$ holds for all $\theta \in \mathbb{R}^p$. Fix a point θ_0 and a number $R > 0$. For any θ satisfying $\|\theta - \theta_0\| \leq R$, we have $\|\mathcal{J}(\theta) - \mathcal{J}(\theta_0)\| \leq \frac{\alpha}{3}$.

Assumption 3 We assume for all $\theta \in \mathbb{R}^p$ obeying $\|\theta - \theta_0\| \leq R$, we have $\|\mathcal{J}(\theta)\| \leq \beta$.

With these assumptions in place we are now ready to state the following result from [Oymak & Soltanolkotabi \(2020\)](#):

Theorem A.11 Given $\theta_0 \in \mathbb{R}^p$, suppose assumptions 1, 2, and 3 hold with

$$R = \frac{3\|f(\theta_0) - y\|_{\ell_2}}{\alpha}. \tag{37}$$

Then, using a learning rate $\eta \leq \frac{1}{3\beta^2}$, all gradient descent updates obey

$$\|f(\theta_\tau) - y\|_{\ell_2} \leq (1 - \eta\alpha^2)^\tau \|f(\theta_0) - y\|_{\ell_2}. \quad (38)$$

We are going to apply this theorem in our case where the parameter is \mathbf{W} , the nonlinear mapping is $f(\mathbf{W}) = \frac{1}{n} \sum_{i=1}^n \mathbf{V} \phi(\mathbf{W} \mathbf{z}_i)$ with $\phi = \text{ReLU}$, and the norm $\|\cdot\|$ set to the operator norm.

Similar to previous part, by using Lemma A.7 we conclude that with probability at least $1 - 3e^{-\frac{m}{72}} - 4k \cdot e^{-c \cdot n} - 2k \cdot e^{-\frac{d}{216}}$, assumption 1 is satisfied with

$$2\alpha := \sigma_v \sigma_z \left(\frac{1}{2} \sqrt{k-9} - 2\sqrt{m} \right) \sqrt{\frac{d + \frac{n-1}{\pi}}{2n}}.$$

Next we show that assumption 2 is valid for α as defined in the previous line and for radius \tilde{R} defined later. First we note that

$$\frac{\alpha}{3} \geq c \cdot \sigma_v \sigma_z \cdot \sqrt{k},$$

where the inequality holds by assuming $K \geq C \cdot m$ for a sufficiently large positive constant C . Now by using equation 35 assumption 2 holds if

$$\frac{5}{4} \sigma_v \sigma_z \sqrt{d} \left(2\sqrt{m} + \sqrt{6 \left(\frac{2kR}{\sigma_w} \right)^{\frac{2}{3}} \log \left(\frac{k}{3 \left(\frac{2kR}{\sigma_w} \right)^{\frac{2}{3}}} \right)} \right) \leq c \cdot \sigma_v \sigma_z \cdot \sqrt{k},$$

which is equivalent to

$$C\sqrt{md} + C\sqrt{d} \cdot \sqrt{6 \left(\frac{2kR}{\sigma_w} \right)^{\frac{2}{3}} \log \left(\frac{k}{3 \left(\frac{2kR}{\sigma_w} \right)^{\frac{2}{3}}} \right)} \leq \sqrt{k}.$$

The first term in the L.H.S. of the inequality above is upper bounded by $\frac{1}{2}\sqrt{k}$ if $k \geq C \cdot md$. For upper bounding the second term it is sufficient to show that

$$C\sqrt{d} \sqrt{6 \left(\frac{2kR}{\sigma_w} \right)^{\frac{2}{3}} \log \left(\frac{k}{3 \cdot \left(\frac{2kR}{\sigma_w} \right)^{\frac{2}{3}}} \right)} \leq \sqrt{k},$$

which by defining $x = 3d \left(\frac{2R}{\sigma_w \sqrt{k}} \right)^{\frac{2}{3}}$ is equivalent to $x \cdot \log \left(\frac{d}{x} \right) \leq C$. Now this last inequality holds if we have $x \leq \frac{c}{\log(d)}$ for a sufficiently small constant c , which by rearranging terms amounts to showing that $R \leq c \cdot \frac{\sigma_w \sqrt{k}}{(d \cdot \log(d))^{\frac{3}{2}}}$. Hence up to this point, we have shown that assumption 2 holds with radius $\tilde{R} := c \cdot \frac{\sigma_w \sqrt{k}}{(d \cdot \log(d))^{\frac{3}{2}}}$, and this implies that it holds for all values of R less than \tilde{R} .

Therefore, we work with the definition of R in equation 37 to show that $R \leq \tilde{R}$ as follows:

$$\begin{aligned} R &= \frac{3 \|f(\theta_0) - y\|_{\ell_2}}{\alpha} \\ &\stackrel{(i)}{\leq} \frac{3}{\alpha} \left(\frac{2}{3} \sigma_v \sigma_w \sigma_z \sqrt{k \cdot d \cdot m} + \|\bar{\mathbf{x}}\|_{\ell_2} \right) \\ &= \frac{2 \left(2\sigma_v \sigma_w \sigma_z \sqrt{k \cdot m \cdot d} + 3 \|\bar{\mathbf{x}}\|_{\ell_2} \right)}{\sigma_v \sigma_z \left(\frac{1}{2} \sqrt{k-9} - 2\sqrt{m} \right) \sqrt{\frac{d + \frac{n-1}{\pi}}{2n}}}, \end{aligned}$$

where in (i) we used Lemma A.10 with $\delta = \frac{1}{3}$. Hence for showing $R \leq \tilde{R}$ it suffices to show that

$$\frac{2 \left(2\sigma_v \sigma_w \sigma_z \sqrt{k \cdot m \cdot d} + 3 \|\bar{\mathbf{x}}\|_{\ell_2} \right)}{\sigma_v \sigma_z \left(\frac{1}{2} \sqrt{k-9} - 2\sqrt{m} \right) \sqrt{\frac{d + \frac{n-1}{\pi}}{2n}}} \leq c \cdot \frac{\sigma_w \sqrt{k}}{(d \cdot \log(d))^{\frac{3}{2}}},$$

which by assuming $k \geq C \cdot m$ simplifies to

$$\sigma_v \sigma_w \sigma_z (d \cdot \log(d))^{\frac{3}{2}} \sqrt{k \cdot m \cdot d} + (d \cdot \log(d))^{\frac{3}{2}} \|\bar{\mathbf{x}}\|_{\ell_2} \leq C \cdot k \cdot \sigma_v \sigma_w \sigma_z.$$

Now this last inequality holds if $k \geq C \cdot m d^4 \log(d)^3$ and by setting $\sigma_v \sigma_w \sigma_z \geq \frac{\|\bar{\mathbf{x}}\|_{\ell_2}}{m d^{\frac{5}{2}} \log d^{\frac{3}{2}}}$.

Therefore Assumption 2 holds for radius R defined in equation 37 with probability at least $1 - n e^{-\frac{m}{2}} - n e^{-c \cdot m d^3 \log(d)^2} - k \cdot e^{-c \cdot n} - e^{-\frac{m}{1500}} - e^{-\frac{k d}{162}}$.

Finally to show assumption 3 holds, we note that for all \mathbf{W} satisfying $\|\mathbf{W} - \mathbf{W}_0\| \leq R$, where the value of R is defined in equation 37, it holds that

$$\begin{aligned} \|\mathcal{J}(\mathbf{W})\| &\leq \|\mathcal{J}(\mathbf{W}_0)\| + \|\mathcal{J}(\mathbf{W}) - \mathcal{J}(\mathbf{W}_0)\| \\ &\leq \|\mathcal{J}(\mathbf{W}_0)\| + \frac{\alpha}{3} \\ &\leq \|\mathcal{J}(\mathbf{W}_0)\| + \frac{\sigma_{\min}(\mathcal{J}(\mathbf{W}_0))}{6} \\ &\leq 2 \|\mathcal{J}(\mathbf{W}_0)\| \\ &\leq 3 \sigma_v \sigma_z \left(\sqrt{k} + 2\sqrt{m} \right) \sqrt{\frac{d + \frac{n-1}{\pi}}{2n}}, \end{aligned}$$

where the last inequality holds by using lemma A.8, hence establishing that assumption 3 holds with

$$\beta = 3 \sigma_v \sigma_z \left(\sqrt{k} + 2\sqrt{m} \right) \sqrt{\frac{d + \frac{n-1}{\pi}}{2n}}$$

with probability at least $1 - e^{-\frac{m}{2}} - 2k \cdot e^{-c \cdot n} - k \cdot e^{-\frac{d}{216}}$, finishing the proof of Theorem 2.2.

B PROOFS OF THE AUXILIARY LEMMAS

In this section, we first provide a proof of Lemma A.6 and next go over the proofs of the key lemmas stated in Section A.5.

B.1 PROOF OF LEMMA A.6

Recall that

$$\begin{aligned} \mathcal{J}(\mathbf{W}) \mathcal{J}(\mathbf{W})^T &= \frac{1}{n^2} \sum_{i,j=1}^n (\mathbf{V} \text{diag}(\phi'(\mathbf{W} \mathbf{z}_i)) \text{diag}(\phi'(\mathbf{W} \mathbf{z}_j) \mathbf{V}^T) (\mathbf{z}_i^T \mathbf{z}_j)) \\ &= \frac{1}{n^2} \mathbf{V} \text{diag}_{\ell=1,\dots,k} \left(\left\| \sum_{i=1}^n \mathbf{z}_i \phi'(\mathbf{w}_\ell^T \mathbf{z}_i) \right\|_{\ell_2}^2 \right) \mathbf{V}^T = \frac{1}{n^2} \mathbf{V} \cdot \mathbf{D}^2 \cdot \mathbf{V}^T, \end{aligned}$$

where \mathbf{D} is a diagonal matrix with entries

$$D_{\ell\ell} = \left\| \sum_{i=1}^n \mathbf{z}_i \phi'(\mathbf{w}_\ell^T \mathbf{z}_i) \right\|_{\ell_2} = \|\mathbf{Z}^T \phi'(\mathbf{Z} \mathbf{w}_\ell)\|_{\ell_2}. \quad (39)$$

The matrix $\mathbf{Z} \in \mathbb{R}^{n \times d}$ contains the \mathbf{z}_i 's in its rows. In order to proceed we are going to analyze the entries of the diagonal matrix \mathbf{D}^2 . We observe that

$$\|\mathbf{Z}^T \phi'(\mathbf{Z} \mathbf{w})\|_{\ell_2}^2 = \underbrace{\left\| \left(\mathbf{I} - \frac{\mathbf{w} \mathbf{w}^T}{\|\mathbf{w}\|^2} \right) \mathbf{Z}^T \phi'(\mathbf{Z} \mathbf{w}) \right\|_{\ell_2}^2}_A + \underbrace{\left\| \frac{\mathbf{w} \mathbf{w}^T}{\|\mathbf{w}\|^2} \mathbf{Z}^T \phi'(\mathbf{Z} \mathbf{w}) \right\|_{\ell_2}^2}_B.$$

First, we compute the expectation of A . We observe that

$$A = \left\| \sum_{i=1}^n \left(\mathbf{I} - \frac{\mathbf{w} \mathbf{w}^T}{\|\mathbf{w}\|^2} \right) \mathbf{z}_i \phi'(\mathbf{w}^T \mathbf{z}_i) \right\|_{\ell_2}^2.$$

Conditioned on \mathbf{w} , $\left(\mathbf{I} - \frac{\mathbf{w}\mathbf{w}^T}{\|\mathbf{w}\|^2}\right) \mathbf{z}_i$ is distributed as $\mathcal{N}\left(0, \sigma_z^2 \left(\mathbf{I} - \frac{\mathbf{w}\mathbf{w}^T}{\|\mathbf{w}\|^2}\right)\right)$ and $\mathbf{w}^T \mathbf{z}_i$ has distribution $\mathcal{N}\left(0, \sigma_z^2 \|\mathbf{w}\|^2\right)$. Moreover, these two random variables are independent, because \mathbf{w} is in the null space of $\mathbf{I} - \frac{\mathbf{w}\mathbf{w}^T}{\|\mathbf{w}\|^2}$. This observation yields

$$\begin{aligned} \mathbb{E}(A) &= \mathbb{E} \left\| \sum_{i=1}^n \left(\mathbf{I} - \frac{\mathbf{w}\mathbf{w}^T}{\|\mathbf{w}\|^2}\right) \mathbf{z}_i \phi'(\mathbf{w}^T \mathbf{z}_i) \right\|_{\ell_2}^2 \\ &= \sum_{i=1}^n \sum_{j=1}^n \mathbb{E} \left(\left\langle \left(\mathbf{I} - \frac{\mathbf{w}\mathbf{w}^T}{\|\mathbf{w}\|^2}\right) \mathbf{z}_i, \left(\mathbf{I} - \frac{\mathbf{w}\mathbf{w}^T}{\|\mathbf{w}\|^2}\right) \mathbf{z}_j \right\rangle \phi'(\mathbf{w}^T \mathbf{z}_i) \phi'(\mathbf{w}^T \mathbf{z}_j) \right) \\ &= \sum_{i=1}^n \mathbb{E} \left(\left\| \left(\mathbf{I} - \frac{\mathbf{w}\mathbf{w}^T}{\|\mathbf{w}\|^2}\right) \mathbf{z}_i \right\|_{\ell_2}^2 \right) \mathbb{E} \left([\phi'(\mathbf{w}^T \mathbf{z}_i)]^2 \right) \\ &= \sum_{i=1}^n \frac{1}{2} (d-1) \sigma_z^2 = \frac{n}{2} (d-1) \sigma_z^2. \end{aligned}$$

Next we show that A is concentrated around its mean. Because $\left(\mathbf{I} - \frac{\mathbf{w}\mathbf{w}^T}{\|\mathbf{w}\|^2}\right) \mathbf{z}_i$ is independent from $\mathbf{w}^T \mathbf{z}_i$, we use \mathbf{z}'_i as an independent copy of \mathbf{z}_i . Hence we can write

$$\begin{aligned} A &= \left\| \left(\mathbf{I} - \frac{\mathbf{w}\mathbf{w}^T}{\|\mathbf{w}\|^2}\right) \mathbf{Z}^T \phi'(\mathbf{Z}\mathbf{w}) \right\|_{\ell_2}^2 \\ &= \left\| \left(\mathbf{I} - \frac{\mathbf{w}\mathbf{w}^T}{\|\mathbf{w}\|^2}\right) \mathbf{Z}^T \phi'(\mathbf{Z}'\mathbf{w}) \right\|_{\ell_2}^2 \\ &= \left\| \sum_{i=1}^n \left(\mathbf{I} - \frac{\mathbf{w}\mathbf{w}^T}{\|\mathbf{w}\|^2}\right) \mathbf{z}_i \phi'(\mathbf{w}^T \mathbf{z}'_i) \right\|_{\ell_2}^2 = \left\| \sum_{i=1}^n \mathbf{g}_i u_i \right\|_{\ell_2}^2, \end{aligned}$$

where $\mathbf{g}_i = \left(\mathbf{I} - \frac{\mathbf{w}\mathbf{w}^T}{\|\mathbf{w}\|^2}\right) \mathbf{z}_i \sim \mathcal{N}\left(0, \sigma_z^2 \left(\mathbf{I} - \frac{\mathbf{w}\mathbf{w}^T}{\|\mathbf{w}\|^2}\right)\right)$ and $u_i = \phi'(\mathbf{w}^T \mathbf{z}'_i) \sim \text{bern}\left(\frac{1}{2}\right)$,⁴ and these are all independent from each other. Note that $\left\| \sum_{i=1}^n \mathbf{g}_i u_i \right\|_{\ell_2}^2$ has the same distribution as $\|\mathbf{g}\|_{\ell_2}^2 \cdot \|\mathbf{u}\|_{\ell_2}^2$, where $\mathbf{g} \sim \mathcal{N}\left(0, \sigma_z^2 \left(\mathbf{I} - \frac{\mathbf{w}\mathbf{w}^T}{\|\mathbf{w}\|^2}\right)\right)$ and \mathbf{u} is a vector with entries u_i . Note that for the norm of \mathbf{u} , the event

$$\frac{n}{2} (1 - \delta) \leq \|\mathbf{u}\|_{\ell_2}^2 \leq \frac{n}{2} (1 + \delta)$$

holds with probability at least $1 - 2e^{-\frac{n\delta^2}{2}}$. Recall that for $\mathbf{g} \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_d)$ and $0 < \delta \leq \frac{1}{2}$ we have

$$\begin{aligned} \mathbb{P} \left(\|\mathbf{g}\|_{\ell_2}^2 \geq (1 + \delta) \mathbb{E} \left(\|\mathbf{g}\|_{\ell_2}^2 \right) \right) &\leq e^{-\frac{d\delta^2}{6}}, \\ \mathbb{P} \left(\|\mathbf{g}\|_{\ell_2}^2 \leq (1 - \delta) \mathbb{E} \left(\|\mathbf{g}\|_{\ell_2}^2 \right) \right) &\leq e^{-\frac{d\delta^2}{4}}. \end{aligned} \tag{40}$$

By applying the union bound and noting that $\mathbb{E} \left(\|\mathbf{g}\|_{\ell_2}^2 \right) = (d-1) \sigma_z^2$, for $0 < \delta \leq \frac{3}{2}$, we obtain that the event

$$\left| \|\mathbf{g}\|_{\ell_2}^2 \|\mathbf{u}\|_{\ell_2}^2 - \frac{n}{2} (d-1) \sigma_z^2 \right| \leq \delta \frac{n}{2} (d-1) \sigma_z^2$$

⁴Here, $\text{bern}\left(\frac{1}{2}\right)$ means that the random variable takes values 0 and 1 each with probability 1/2.

holds with probability at least $1 - 2e^{-\frac{n\delta^2}{18}} - 2e^{-\frac{d\delta^2}{54}}$.

In order to analyze B , we first note that

$$\begin{aligned} B &= \left\| \frac{\mathbf{w}\mathbf{w}^T}{\|\mathbf{w}\|^2} \mathbf{Z}^T \phi'(\mathbf{Z}\mathbf{w}) \right\|_{\ell_2}^2 = \left| \frac{\mathbf{w}^T}{\|\mathbf{w}\|} \mathbf{Z}^T \phi'(\mathbf{Z}\mathbf{w}) \right|^2 \\ &= \left| \left\langle \mathbf{Z} \frac{\mathbf{w}}{\|\mathbf{w}\|}, \phi'(\mathbf{Z}\mathbf{w}) \right\rangle \right|^2 \\ &= \left| \langle \mathbf{g}, \phi'(\|\mathbf{w}\|\mathbf{g}) \rangle \right|^2 \\ &= \left| \sum_{i=1}^n \mathbf{g}_i \cdot \mathbb{1}_{(\mathbf{g}_i \geq 0)} \right|^2 = \left(\sum_{i=1}^n \text{ReLU}(\mathbf{g}_i) \right)^2, \end{aligned}$$

where $\mathbf{g}_i = z_i^T \frac{\mathbf{w}}{\|\mathbf{w}\|} \sim \mathcal{N}(0, \sigma_z^2)$. It follows that

$$\begin{aligned} \mathbb{E}(B) &= \mathbb{E} \left(\sum_{i=1}^n \text{ReLU}(\mathbf{g}_i) \right)^2 \\ &= \sum_{i=1}^n \mathbb{E}(\text{ReLU}^2(\mathbf{g}_i)) + \sum_{i \neq j} \mathbb{E}(\text{ReLU}(\mathbf{g}_i) \text{ReLU}(\mathbf{g}_j)) \\ &= \sigma_z^2 \left(\frac{n}{2} + \frac{n(n-1)}{2\pi} \right), \end{aligned}$$

which results in

$$\mathbb{E}(\mathbf{D}_{\ell\ell}^2) = \mathbb{E}(A) + \mathbb{E}(B) = \sigma_z^2 \left(\frac{nd}{2} + \frac{n(n-1)}{2\pi} \right), \quad 1 \leq \ell \leq k.$$

Next, in order to show that B concentrates around its mean, we note that $\text{ReLU}(\mathbf{g}_i)$ is a sub-Gaussian random variable with ψ_2 -norm $C\sigma_z$, where C is a fixed constant. Therefore $X = \sum_{i=1}^n \text{ReLU}(\mathbf{g}_i)$ is sub-Gaussian with ψ_2 -norm $C\sqrt{n}\sigma_z$. By the sub-exponential tail bound for $X^2 - \mathbb{E}(X^2)$ we obtain

$$\mathbb{P}(|B - \mathbb{E}(B)| \geq t) \leq 2e^{-c\frac{t}{n\sigma_z^2}}.$$

Finally by putting these results together and using union bounds we have

$$\mathbb{P}\{|\mathbf{D}_{\ell\ell}^2 - \mathbb{E}(\mathbf{D}_{\ell\ell}^2)| \geq \delta \mathbb{E}(\mathbf{D}_{\ell\ell}^2)\} \leq 2e^{-\frac{n\delta^2}{18}} + 2e^{-\frac{d\delta^2}{54}} + 2e^{-c_1 n\delta}, \quad 0 \leq \delta \leq \frac{3}{2},$$

finishing the proof of Lemma A.6.

B.2 PROOF OF LEMMA A.7

Our main tool for bounding the minimum singular value of the Jacobian mapping will be the following lemma from [Soltanolkotabi \(2019\)](#):

Lemma B.1 *Let $\mathbf{d} \in \mathbb{R}^k$ be a fixed vector with nonzero entries and $\mathbf{D} = \text{diag}(\mathbf{d})$. Also, let $\mathbf{A} \in \mathbb{R}^{k \times m}$ have i.i.d. $\mathcal{N}(0, 1)$ entries and $\mathcal{T} \subseteq \mathbb{R}^m$. Define*

$$b_k(\mathbf{d}) = \mathbb{E}[\|\mathbf{D}\mathbf{g}\|_{\ell_2}],$$

where $\mathbf{g} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_k)$. Also define

$$\sigma(\mathcal{T}) := \max_{\mathbf{v} \in \mathcal{T}} \|\mathbf{v}\|_{\ell_2}.$$

Then for all $\mathbf{u} \in \mathcal{T}$ we have

$$\|\mathbf{D}\mathbf{A}\mathbf{u}\|_{\ell_2} - b_k(\mathbf{d})\|\mathbf{u}\|_{\ell_2} \leq \|\mathbf{d}\|_{\ell_\infty} \omega(\mathcal{T}) + \eta$$

with probability at least $1 - 6e^{-\frac{\eta^2}{8\|\mathbf{d}\|_{\ell_\infty}^2 \sigma^2(\mathcal{T})}}$.

In order to apply this lemma, we set the elements of \mathbf{d} to be $D_{\ell\ell}$ as in equation 39 and choose $\mathcal{T} = S^{m-1}$ and $\mathbf{A} = \mathbf{V}^T \in \mathbb{R}^{k \times m}$ with $\mathcal{N}(0, \sigma_v^2)$ entries. It follows that

$$b_k(\mathbf{d}) = \mathbb{E} \|\mathbf{D}\mathbf{g}\|_{\ell_2} = \sqrt{\mathbb{E} \left(\|\mathbf{D}\mathbf{g}\|_{\ell_2}^2 \right) - \text{Var} \left(\|\mathbf{D}\mathbf{g}\|_{\ell_2} \right)},$$

where

$$\mathbb{E} \left(\|\mathbf{D}\mathbf{g}\|_{\ell_2}^2 \right) = \|\mathbf{d}\|_{\ell_2}^2 = \sum_{\ell=1}^k \mathbf{D}_{\ell\ell}^2.$$

We are going to use the fact that for a B -Lipschitz function ϕ and normal random variable $g \sim \mathcal{N}(0, 1)$, based on the Poincare inequality (Ledoux, 2001) we have $\text{Var}(\phi(g)) \leq B^2$. By noting that for a diagonal matrix \mathbf{D}

$$\left| \|\mathbf{D}\mathbf{x}\|_{\ell_2} - \|\mathbf{D}\mathbf{y}\|_{\ell_2} \right| \leq \|\mathbf{D}\mathbf{x} - \mathbf{D}\mathbf{y}\|_{\ell_2} \leq \|\mathbf{d}\|_{\ell_\infty} \|\mathbf{x} - \mathbf{y}\|_{\ell_2},$$

we get

$$\begin{aligned} \mathbb{E} \|\mathbf{D}\mathbf{g}\|_{\ell_2} &= \sqrt{\mathbb{E} \left(\|\mathbf{D}\mathbf{g}\|_{\ell_2}^2 \right) - \text{Var}(\|\mathbf{D}\mathbf{g}\|_{\ell_2})} \\ &\geq \sqrt{\|\mathbf{d}\|_{\ell_2}^2 - \|\mathbf{d}\|_{\ell_\infty}^2}. \end{aligned}$$

This combined with $\omega(S^{m-1}) \leq \sqrt{m}$ and Lemma B.1 yields that the event

$$\sigma_{\min}(\mathbf{V}\mathbf{D}) \geq \sigma_v \left(\sqrt{\|\mathbf{d}\|_{\ell_2}^2 - \|\mathbf{d}\|_{\ell_\infty}^2} - \|\mathbf{d}\|_{\ell_\infty} \sqrt{m} - \eta \right) \quad (41)$$

holds with probability at least $1 - 3e^{-\frac{\eta^2}{8\|\mathbf{d}\|_{\ell_\infty}^2}}$.

Next, using the concentration bound for $\mathbf{D}_{\ell\ell}^2$, which we obtained in Section B.1, we bound $\|\mathbf{d}\|_{\ell_2}^2$ and $\|\mathbf{d}\|_{\ell_\infty}$, where we have set $\mathbf{d}_i = \mathbf{D}_{ii}$ for $1 \leq i \leq k$. For $0 \leq \delta \leq \frac{3}{2}$ we compute that

$$\begin{aligned} \mathbb{P} \left(\max_{1 \leq i \leq k} \mathbf{d}_i \geq (1 + \delta) \sqrt{\mathbb{E}[\mathbf{d}_i^2]} \right) &= \mathbb{P} \left(\bigcup_{i=1}^k \mathbf{d}_i^2 \geq (1 + \delta)^2 \mathbb{E}[\mathbf{d}_i^2] \right) \\ &\leq k \cdot \mathbb{P} \left(\mathbf{d}_i^2 \geq (1 + \delta)^2 \mathbb{E}[\mathbf{d}_i^2] \right) \\ &\leq k \cdot \mathbb{P} \left(\mathbf{d}_i^2 \geq (1 + \delta) \mathbb{E}[\mathbf{d}_i^2] \right) \leq k \cdot \left(e^{-\frac{n\delta^2}{18}} + e^{-\frac{d\delta^2}{54}} + e^{-c_1 n\delta} \right), \end{aligned} \quad (42)$$

as well as

$$\begin{aligned} \mathbb{P} \left(\|\mathbf{d}\|_{\ell_2} \leq (1 - \delta) \sqrt{k} \sqrt{\mathbb{E}[\mathbf{d}_i^2]} \right) &\leq \mathbb{P} \left(\bigcup_{i=1}^k \mathbf{d}_i^2 \leq (1 - \delta)^2 \mathbb{E}[\mathbf{d}_i^2] \right) \\ &\leq k \cdot \mathbb{P} \left(\mathbf{d}_i^2 \leq (1 - \delta)^2 \mathbb{E}[\mathbf{d}_i^2] \right) \\ &\leq k \cdot \mathbb{P} \left(\mathbf{d}_i^2 \leq (1 - \delta) \mathbb{E}[\mathbf{d}_i^2] \right) \leq k \cdot \left(e^{-\frac{n\delta^2}{18}} + e^{-\frac{d\delta^2}{54}} + e^{-c_1 n\delta} \right). \end{aligned} \quad (43)$$

Finally by replacing η with $\eta \|\mathbf{d}\|_{\ell_\infty} \sqrt{m}$ in equation 41, combined with equation 42 and equation 43, for a random \mathbf{W}_0 with i.i.d. $\mathcal{N}(0, \sigma_w^2)$ entries we have:

$$\begin{aligned} \sigma_{\min}(\mathcal{J}(\mathbf{W}_0)) &= \frac{1}{n} \sigma_{\min}(\mathbf{V}\mathbf{D}) \\ &\geq \frac{\sigma_v}{n} \left(\sqrt{(1 - \delta)^2 k - (1 + \delta)^2} - \sqrt{m} (1 + \eta) (1 + \delta) \right) \sqrt{\mathbb{E}[\mathbf{d}_i^2]} \\ &= \left(\sqrt{(1 - \delta)^2 k - (1 + \delta)^2} - \sqrt{m} (1 + \eta) (1 + \delta) \right) \sigma_v \sigma_z \sqrt{\frac{d + \frac{n-1}{\pi}}{2n}}, \quad 0 \leq \delta \leq \frac{3}{2}, \end{aligned}$$

with probability at least $1 - 3e^{-\frac{\eta^2 m}{8}} - 2k \cdot \left(e^{-\frac{n\delta^2}{18}} + e^{-\frac{d\delta^2}{54}} + e^{-c_1 n\delta} \right)$. This completes the proof of Lemma A.7.

B.3 PROOF OF LEMMA A.8

Recall that

$$\mathcal{J}(\mathbf{W})\mathcal{J}(\mathbf{W})^T = \frac{1}{n^2} \mathbf{V} \operatorname{diag}_{\ell=1, \dots, k} \left(\left\| \sum_{i=1}^n \mathbf{z}_i \phi'(\mathbf{w}_\ell^T \mathbf{z}_i) \right\|_{\ell_2}^2 \right) \mathbf{V}^T = \frac{1}{n^2} \mathbf{V} \cdot \mathbf{D}^2 \cdot \mathbf{V}^T,$$

which implies that

$$\|\mathcal{J}(\mathbf{W}_0)\| = \frac{1}{n} \|\mathbf{V} \cdot \mathbf{D}\| \leq \frac{1}{n} \|\mathbf{V}\| \|\mathbf{D}\|.$$

For matrix $\mathbf{V} \in \mathbb{R}^{m \times k}$ with i.i.d $\mathcal{N}(0, \sigma_v^2)$ the event

$$\|\mathbf{V}\| \leq \sigma_v \left(\sqrt{k} + 2\sqrt{m} \right)$$

holds with probability at least $1 - e^{-\frac{m}{2}}$. Regarding matrix \mathbf{D} by repeating equation 42 the following event

$$\|\mathbf{D}\| = \max_{1 \leq i \leq k} D_{ii} \leq (1 + \delta) \sqrt{\mathbb{E}[D_{ii}^2]} = (1 + \delta) \sigma_z \sqrt{\frac{nd}{2} + \frac{n(n-1)}{2\pi}}, \quad 0 \leq \delta \leq \frac{3}{2}$$

holds with probability at least $1 - k \cdot \left(e^{-\frac{n\delta^2}{18}} + e^{-\frac{d\delta^2}{54}} + e^{-c_1 n \delta} \right)$. Putting these together it yields that the event

$$\|\mathcal{J}(\mathbf{W}_0)\| \leq (1 + \delta) \sigma_v \sigma_z \left(\sqrt{k} + 2\sqrt{m} \right) \sqrt{\frac{d + \frac{n-1}{\pi}}{2n}}, \quad 0 \leq \delta \leq \frac{3}{2}$$

holds with probability at least $1 - e^{-\frac{m}{2}} - k \cdot \left(e^{-\frac{n\delta^2}{18}} + e^{-\frac{d\delta^2}{54}} + e^{-c_1 n \delta} \right)$, finishing the proof of Lemma A.8.

B.4 PROOF OF LEMMA A.10

First, note that if \mathbf{W} has i.i.d. $\mathcal{N}(0, \sigma_w^2)$ entries and $\mathbf{V}, \mathbf{W}, \mathbf{Z}$ are all independent, then $\|f(\mathbf{W})\|_{\ell_2} = \frac{1}{n} \|\mathbf{V} \phi(\mathbf{W} \mathbf{Z}^T) \mathbf{1}_{n \times 1}\|_{\ell_2}$ has the same distribution as $\|\mathbf{v}\|_{\ell_2} \|\mathbf{a}\|_{\ell_2}$, where $\mathbf{v} \sim \mathcal{N}(0, \sigma_v^2 \mathbf{I}_m)$ and $\mathbf{a} = \frac{1}{n} \phi(\mathbf{W} \mathbf{Z}^T) \mathbf{1}$ has independent sub-Gaussian entries, so its ℓ_2 -norm is concentrated. Note that conditioned on \mathbf{W} , $a_i = \frac{1}{n} \sum_{j=1}^n \operatorname{ReLU}(\mathbf{z}_j^T \mathbf{w}_i)$ is sub-Gaussian with $\|a_i\|_{\psi_2} = C \frac{\|\mathbf{w}_i\|_{\ell_2} \sigma_z}{\sqrt{n}}$, and it is concentrated around $\mathbb{E} a_i = \frac{1}{\sqrt{2\pi}} \|\mathbf{w}_i\|_{\ell_2} \sigma_z$. This gives

$$\mathbb{P} \{ a_i \leq (1 + \delta) \mathbb{E} a_i \} \geq 1 - e^{-c \frac{\delta^2 (\mathbb{E} a_i)^2}{\|a_i\|_{\psi_2}^2}} = 1 - e^{-cn\delta^2},$$

which implies that

$$\mathbb{P} \{ a_i^2 \leq (1 + 3\delta) (\mathbb{E} a_i)^2 \} \geq \mathbb{P} \{ a_i^2 \leq (1 + \delta)^2 (\mathbb{E} a_i)^2 \} \geq 1 - e^{-cn\delta^2}, \quad 0 \leq \delta \leq 1.$$

Due to the union bound we get that

$$\begin{aligned} \mathbb{P} \left\{ \|\mathbf{a}\|_{\ell_2}^2 \geq (1 + \delta) \sum_{i=1}^k (\mathbb{E} a_i)^2 \right\} &\leq \mathbb{P} \left\{ \bigcup_{i=1}^k a_i^2 \geq (1 + \delta) (\mathbb{E} a_i)^2 \right\} \\ &\leq \sum_{i=1}^k \mathbb{P} \{ a_i^2 \geq (1 + \delta) (\mathbb{E} a_i)^2 \} \leq k \cdot e^{-cn(\delta/3)^2}, \quad 0 \leq \delta \leq 3. \end{aligned}$$

By substituting $\sum_{i=1}^k (\mathbb{E} a_i)^2 = \frac{1}{2\pi} \sigma_z^2 \|\mathbf{W}\|_F^2$ this shows

$$\mathbb{P} \left\{ \|\mathbf{a}\|_{\ell_2} \leq (1 + \delta) \frac{1}{\sqrt{2\pi}} \sigma_z \|\mathbf{W}\|_F \right\} \geq \mathbb{P} \left\{ \|\mathbf{a}\|_{\ell_2}^2 \leq (1 + \delta) \frac{1}{2\pi} \sigma_z^2 \|\mathbf{W}\|_F^2 \right\} \geq 1 - k \cdot e^{-cn(\delta/3)^2}, \quad 0 \leq \delta \leq 3.$$

We also have the following result for $\mathbf{v} \sim \mathcal{N}(0, \sigma_v^2 \mathbf{I}_m)$

$$\mathbb{P} \left\{ \|\mathbf{v}\|_{\ell_2} \leq (1 + \delta) \sigma_v \sqrt{m} \right\} \geq 1 - e^{-\frac{\delta^2 m}{2}}.$$

By combining the above results we obtain

$$\begin{aligned} \mathbb{P} \left\{ \|\mathbf{a}\|_{\ell_2} \|\mathbf{v}\|_{\ell_2} \leq (1 + \delta) \frac{1}{\sqrt{2\pi}} \sigma_v \sigma_z \sqrt{m} \|\mathbf{W}\|_F \right\} &\geq \mathbb{P} \left\{ \|\mathbf{a}\|_{\ell_2} \|\mathbf{v}\|_{\ell_2} \leq (1 + \delta/3)^2 \frac{1}{\sqrt{2\pi}} \sigma_v \sigma_z \sqrt{m} \|\mathbf{W}\|_F \right\} \\ &\geq 1 - k \cdot e^{-cn(\delta/9)^2} - e^{-\frac{(\delta/3)^2 m}{2}}, \quad 0 \leq \delta \leq 3. \end{aligned}$$

Furthermore, we can bound $\|\mathbf{W}\|_F$ by the tail inequality

$$\mathbb{P} \left\{ \|\mathbf{W}\|_F \leq (1 + \delta) \sigma_w \sqrt{kd} \right\} \geq 1 - e^{-\frac{\delta^2 kd}{2}}.$$

Hence, by combining the last two results we have that

$$\begin{aligned} \mathbb{P} \left\{ \|\mathbf{a}\|_{\ell_2} \|\mathbf{v}\|_{\ell_2} \leq (1 + \delta) \frac{1}{\sqrt{2\pi}} \sigma_v \sigma_z \sigma_w \sqrt{k \cdot d \cdot m} \right\} &\geq \mathbb{P} \left\{ \|\mathbf{a}\|_{\ell_2} \|\mathbf{v}\|_{\ell_2} \leq (1 + \delta/3)^2 \frac{1}{\sqrt{2\pi}} \sigma_v \sigma_z \sigma_w \sqrt{k \cdot d \cdot m} \right\} \\ &\geq 1 - k \cdot e^{-cn(\delta/27)^2} - e^{-\frac{(\delta/9)^2 m}{2}} - e^{-\frac{(\delta/3)^2 kd}{2}}, \quad 0 \leq \delta \leq 3. \end{aligned}$$

Therefore, due to the triangle inequality the event

$$\|f(\mathbf{W}_0) - \bar{\mathbf{x}}\|_{\ell_2} \leq (1 + \delta) \frac{1}{\sqrt{2\pi}} \sigma_v \sigma_w \sigma_z \sqrt{k \cdot d \cdot m} + \|\bar{\mathbf{x}}\|_{\ell_2}, \quad 0 \leq \delta \leq 3$$

holds with probability at least $1 - k \cdot e^{-c_2 n(\delta/27)^2} - e^{-\frac{(\delta/9)^2 m}{2}} - e^{-\frac{(\delta/3)^2 kd}{2}}$ for some positive constant c_2 , completing the proof of Lemma A.10.

C ADDITIONAL EXPERIMENTS

Effect of single component overparameterization: In Section 3 of the main paper, we performed experiments in the setting where the size of generator and discriminator are held roughly the same (both discriminator and generator uses the same value of k). In this part, we analyze single-component overparameterization where we study the effect of overparameterization when one of the components (generator / discriminator) has varying k , while the other component uses the standard value of k (64 for DCGAN and 128 for Resnet GAN). The FID variation of single-component overparameterization are shown in Fig. 7. We observe similar trends as the previous case where increasing overparameterization leads to improved FID scores. Interestingly, increasing the value of k beyond the default value used in the other component leads to a slight drop in performance. Hence, choosing comparable sizes of discriminator and generator models is recommended.

D EXPERIMENTAL DETAILS

The model architectures we use this in the experiments are shown in Figure 8. In both DCGAN and Resnet-based GANs, the parameter k controls the number of convolutional filters in each layer. The larger the value of k is, the more overparameterized the models are.

Optimization: Both DCGAN and Resnet-based GAN models are optimized using the commonly used hyper-parameters: Adam with learning rate 0.0001 and betas (0.5, 0.999) for DCGAN, gradient penalty of 10 and 5 critic iterations per generator’s iteration for both DCGAN and Resnet-based GAN models. Models are trained for 300,000 iterations with a batch size of 64.

E NEAREST NEIGHBOR VISUALIZATION

In this section, we visualize the nearest neighbors of samples generated using GAN models trained with different levels of overparameterization. More specifically, we trained a DCGAN model with $k = 8$ and $k = 128$, synthesize random samples from the trained model and query the nearest neighbors in the training set. The plot of obtained samples is shown in Figure. 10. We observe that overparameterized models generate samples with high diversity.

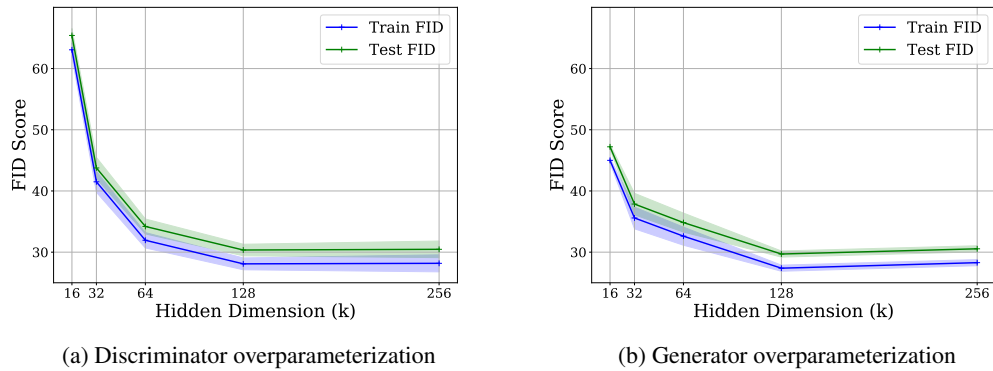


Figure 7: **Single Component Overparameterization Results:** We plot FID scores of Resnet GAN as the hidden dimension of one of the components are varied, while the hidden dimension of other component is held fixed. Even in this case, overparameterization improves model convergence.

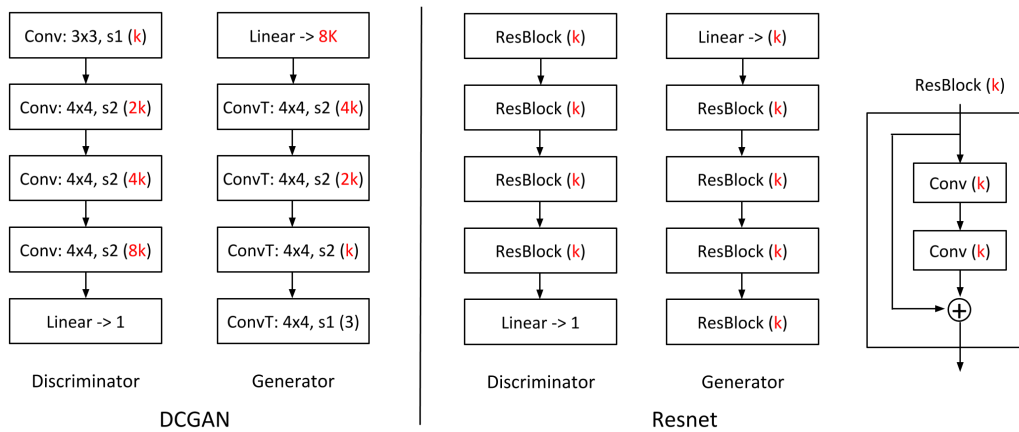


Figure 8: **Architectures used in over-parameterization experiments.** The number of out-channels in convolutional layers is indicated in red. Parameter k controls the width of the architectures – larger the k , more over-parameterized the models are.

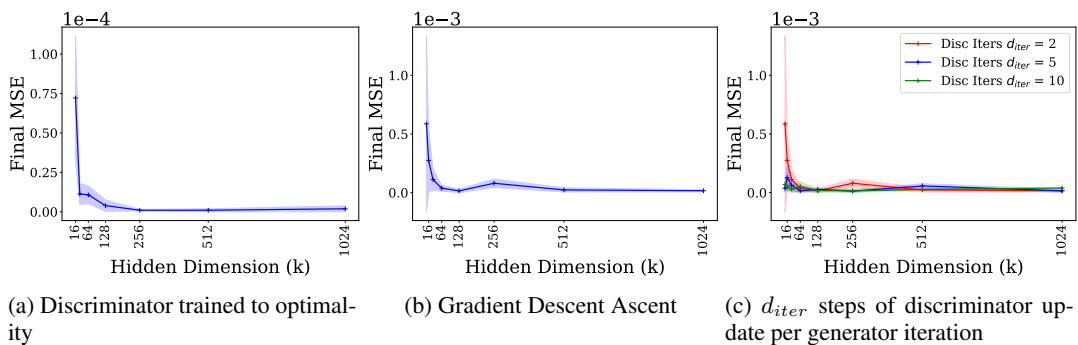


Figure 9: **Convergence plot** GAN model trained on the Two-Moons dataset, with linear discriminator and 1-hidden layer generator as the hidden dimension (k) increases. Over-parameterized models show improved convergence

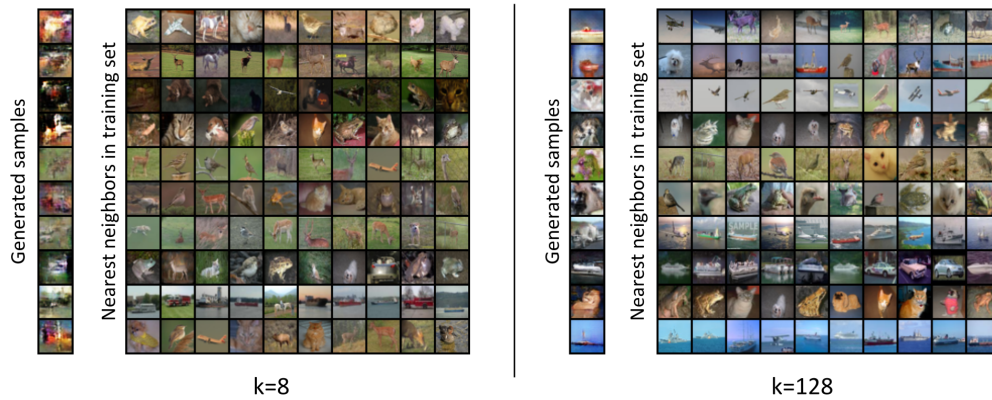


Figure 10: **Nearest neighbor visualization.** We visualize the nearest neighbor samples in training set for generations from DCGAN model trained on CIFAR-10 dataset. Left panel shows DCGAN trained with $k = 8$, while the right one shows the one with $k = 128$. We observe that overparameterized models generate samples with high diversity.